

UNIVERZA V LJUBLJANI
FAKULTETA ZA RAČUNALNIŠTVO IN INFORMATIKO

Anže Brvar

Elektronsko trgovanje na valutnem trgu s pomočjo Twitterja

MAGISTRSKO DELO
MAGISTRSKI PROGRAM DRUGE STOPNJE
RAČUNALNIŠTVO IN INFORMATIKA

MENTOR: izr. prof. dr. Polona Oblak

SOMENTOR: prof. dr. Blaž Zupan

Ljubljana, 2015

AVTORSKE PRAVICE. Rezultati magistrskega dela so intelektualna lastnina avtorja in Fakultete za računalništvo in informatiko Univerze v Ljubljani. Za objavljane ali izkoriščanje rezultatov magistrskega dela je potrebno pisno soglasje avtorja, Fakultete za računalništvo in informatiko ter mentorjev.

©2015 ANŽE BRVAR

IZJAVA O AVTORSTVU MAGISTRSKEGA DELA

Spodaj podpisani Anže Brvar sem avtor magistrskega dela z naslovom:

Elektronsko trgovanje na valutnem trgu s pomočjo Twitterja

S svojim podpisom zagotavljam, da:

- sem magistrsko delo izdelal samostojno pod mentorstvom izr. prof. dr. Polone Oblak in somentorstvom prof. dr. Blaža Zupana,
- so elektronska oblika magistrskega dela, naslov (slovenski, angleški), povzetek (slovenski, angleški) ter ključne besede (slovenske, angleške) identični s tiskano obliko magistrskega dela,
- soglašam z javno objavo elektronske oblike magistrskega dela v zbirki "Dela FRI".

V Ljubljani, 6. novembra 2015

Podpis avtorja:

ZAHVALA

*Za mentorstvo, pomoč, podporo in strokovne nasvete se zahvaljujem mentorici
izr. prof. dr. Poloni Oblak in somentorju prof. dr. Blažu Zupanu.*

*Iskreno se zahvaljujem tudi svojim staršem in Lari za podporo in potrpežljivost
skozi vsa leta študija.*

Anže Brvar, 2015

Kazalo

Povzetek	i
Abstract	iii
1 Uvod	1
2 Valutni trg in trgovanje	5
2.1 Opis trgovanja	6
2.2 Tehnični indikatorji	8
3 Podatki	13
3.1 Twitter	13
3.2 Pretekle vrednosti valutnih tečajev	16
4 Uporabljene metode	19
4.1 Predstavitev besedilnih dokumentov	20
4.2 Napovedovanje	28
5 Eksperimentalno vrednotenje uspešnosti napovednih tehnik	37
5.1 Klasifikacijska točnost	38
5.2 Rezultati	39
6 Simulacija trgovanja	43
6.1 Sistem za simulacijo trgovanja	44
6.2 Izbira primerne razreda	48

KAZALO

6.3	Primerjava različnih predstavitev besedila	50
6.4	Izbira ustreznega števila učnih primerov	50
6.5	Izbira ustreznih uteži	53
6.6	Združevanje napovedi	54
6.7	Najboljših 10 napovednih modelov	58
6.8	Primerjava s trgovanjem na valutnem trgu	59
7	Sklepne ugotovitve	63
A	Seznam rezultatov simulacij	73

Povzetek

V magistrskem delu smo raziskovali uspešnost elektronskega trgovanja na valutnem trgu z metodami strojnega učenja. Primerjali smo uspešnost razvitih algoritmov, ki trgujejo s pomočjo objav (tvitov) na Twitterju, in takih, ki za učne podatke uporabijo pretekle vrednosti valutnih tečajev in tehničnih indikatorjev. Za transformacijo besedil v atributni zapis smo poleg znanih metod preizkusili tudi vektorje besed word2vec. Razvite metode transformacije besedil in njihove parametre smo najprej ovrednotili na sorodnem problemu zaznavanja sentimenta tvitov, nato pa jih preizkusili v trgovanju v simulacijskem okolju. Napovedi razvitih metod smo izboljšali z metodami za združevanje napovedi in tako dosegli do 250% vrednost začetnih sredstev pri simulaciji v obdobju zadnjih petih let. V delu poročamo o najprimernejši izbiri parametrov, ki imajo velik vpliv na uspešnost elektronskega trgovanja. Ugotovili smo, da je Twitter bolj primeren vir informacij za uspešno elektronsko trgovanje kot pretekle vrednosti valutnih tečajev.

Ključne besede

valutno trgovanje, forex, twitter, strojno učenje, word2vec, napovedovanje, simulacija

Abstract

In this thesis we study the performance of electronic trading algorithms with a help of machine learning methods. We compare the performance of developed trading algorithms that trade based on posts (tweets) on Twitter with those that trade based on historic foreign exchange values and technical indicators. Besides the well known methods for text transformation to attribute notation we also use word2vec word vectors. We evaluate all the developed text transformation methods and their parameters, first on simpler but related tweet sentiment detection problem and later with trading in simulation environment. We improve developed models' predictions with the prediction combining techniques and we achieve up to 250% of initial funds at simulation in the period of last five years. The results show that Twitter is a better source of trading information than foreign exchange rates and technical indicators.

Keywords

foreign exchange, forex, twitter, machine learning, word2vec, prediction, simulation

Poglavje 1

Uvod

Sprotna analiza spletnih in javnih virov, kot sta Twitter in spletne novice, je danes zelo aktualno področje raziskav. Obstajajo tehnike, s katerimi je moč napovedovati gibanje različnih finančnih instrumentov [1, 2, 3, 4, 5], prodajo avtomobilov [6], izide volitev [7] in celo pojav epidemij [8].

V zadnjem času se pojavlja vedno več literature, ki opisuje bolj ali manj uspešne sisteme za trgovanje na različnih finančnih trgih [1]. Največji finančni trg na svetu je trg valutnega trgovanja. Dolga leta je bilo valutno trgovanje omejeno na nekaj večjih bank, zadnja leta pa je trgovanje vedno bolj dostopno tudi vlagateljem in posameznikom, ki želijo špekulirati in trgovati z valutnimi pari [9].

Z razvojem spletnih tehnologij se je razširila tudi uporaba elektronskega trgovanja [10]. Posredniki valutnega trgovanja so v zadnjih nekaj letih začeli ponujati veliko orodij za preprosto izdelavo samodejnih elektronskih sistemov za trgovanje na podlagi ročno definiranih pravil. Nassirtoussi in sodelavci [1] tako predstavijo svoj sistem, ki na podlagi naslovov finančnih novic napoveduje gibanje valutnega tečaja EURUSD in doseže 83% natančnost napovedovanja smeri gibanja. Avtorji za razliko od večine preizkušajo napovedovanje na krajših zamikih, saj izdelajo napoved o vrednosti tečaja od ene do treh ur po izdaji novice. V svojem članku navedejo 24 sorodnih del, ki prav tako opisujejo poskuse napovedovanja različnih finančnih in borznih

tečajev s pomočjo tehnične ali temeljne analize, ter na enem mestu predstavijo metode, uporabljene pri napovedovanju. Te vključujejo tehnike izbora atributov, tehnike čiščenja besedila, različne predstavitve besedila, metode strojnega učenja in različne zasnove algoritmov za napovedovanje. Avtorji so v svoji implementaciji uporabili standardne tehnike čiščenja besedila, predstavitev z vrečo besed in dodatno izboljšavo s tehniko uteževanja TF-IDF, ki je podrobneje opisana v poglavju 4.1.1.

Slabost njihovega raziskovanja je majhna učna množica. Velikost testne množice, s katero poročajo o 83% natančnosti je zgolj 12 primerov, ocena točnosti njihovega postopka je lahko zato plod naključja.

Nassirtoussi in sodelavci sicer uporabljajo še nekaj tehnik predpriprave besedil. Uporabijo slovarje sopomenk in protipomenk, ter posebno tehniko uteževanja primerov. Za napovedovanje gibanja tečajev so uporabili finančne novice, navajajo pa tudi, da je ena ura po objavi novice dovolj, da se njen vpliv pozna na vrednosti valutnega tečaja. To smo želeli preveriti tudi v našem delu, vendar smo zaradi razlogov, opisanih v poglavju 3.1.3, za vse eksperimente uporabili enodnevni interval.

Papaioannou in sodelavci [11] poskušajo napovedati dnevno spremembo valutnega tečaja EURUSD na podlagi javno objavljenih naročil na družbenem omrežju Twitter. Za model strojnega učenja preizkusijo dve različni linearni metodi in umetne nevronske mreže. Makrehchi in sodelavci [12] gradijo učni model s pomočjo preteklih odmevnih dogodkov, vrednosti tečajev na borzi in kratkih javnih besedil na omrežju Twitter, ki jih imenujemo tviti (angl. *tweet*). Naučeni model nato uporabijo za napovedovanja čustev, kjer so vhodni podatki tviti tistega dne, napoved pa čustvena naravnost. Glede na njihove rezultate je napovedana čustvena naravnost močno korelirana z gibanjem borznih tečajev. Za vhodne podatke so uporabili le tvite, ki so vsebovali kakšno od ključnih besed napovedovanega borznega tečaja. Na primer, za podjetje Apple so uporabili tvite, ki so vsebovali eno od besed *Apple Inc*, *AAPL* in *appl*. Na podlagi tega članka smo dobili idejo za pridobivanje množice tvitov, ki so omejeni na finančni svet.

Kaya in sodelavci v članku [3] prav tako opisujejo napovedovanje gibanja borznih tečajev na podlagi finančnih novic. Članke očistijo veznikov, jih predstavijo s pari besed, ki se pojavljajo v isti novici in izluščijo le nekaj najbolj informativnih parov, ki jih nato uporabijo za klasifikacijo. Napovedni model po njihovih analizah dosega 61% natančnost.

Obstaja veliko literature na temo zaznavanja čustvene naravnosti. Terana in sodelavci [13] jo zaznavajo s pomočjo emotikonov, Godbole [14] opisuje način zaznavanja z vnaprej podanim slovarjem, ki ga nato razširjajo s pomočjo sopomenk in protipomenk, podobno kot Nassirtoussi in sodelavci [1]. Fong in sodelavci [15] opisujejo zaznavanje čustev z umetnimi nevronskimi mrežami, kjer za vsako čustvo izdelajo ločen model.

V magistrskem delu smo se preizkusili v elektronskem trgovanju na valutnem trgu in pri tem uporabili metode strojnega učenja v kombinaciji s tehnično in temeljno analizo. Tehnična analiza (angl. *technical analysis*) je pristop odločanja o trgovanju zgolj na podlagi zgodovinskih vrednosti tečajev, temeljna analiza (angl. *fundamental analysis*) pa je pristop trgovanja, ki temelji na predpostavkah o tečaju, ki jih lahko pridobimo iz tipično nestrukturiranega teksta (poročila bank, vladnih uradov ali druga zapisana mnenja posameznikov) [1]. Predvsem nas je zanimalo, kateri od pristopov je bolj primeren za elektronsko trgovanje. Preizkusili smo tudi združevanje napovedi izdelanih algoritmov in tako združili pristopa tehnične in temeljne analize, ki sta podrobneje opisana v poglavju 2.1.

Omejili smo se na pridobljene informacije o dogajanju na družbenem omrežju Twitter, kjer več kot 300 milijonov uporabnikov objavlja svoje tvite. Za predstavitev besedila oz. tvitov smo uporabili že znane metode (vreča besed, vreča nizov), preizkusili pa smo tudi predstavitev z dokaj novim pristopom word2vec, ki vključuje vektorje besed. Word2vec vektor besede je vektorska predstavitev besede, ki ohranja podobnosti in relacije med besedami, kar je pri atributni predstavitvi besedila zaželjena lastnost. Omenjene predstavitve besedil smo najprej preizkusili na enostavnejšem problemu - napovedovanju sentimenta in se prepričali, da delujejo zadovoljivo.

Implementirane predstavitve besedil v kombinaciji z nekaj znanimi metodami strojnega učenja smo nato uporabili še na našem glavnem problemu - napovedovanju gibanja valutnih tečajev oz. trgovanju z valutnim parom EURUSD. Za potrebe eksperimentov smo izdelali simulacijsko trgovalno okolje (poglavje 6), s katerim posnemamo resnično trgovanje. Pokazali smo, da obstaja razlika med teoretičnim testiranjem uspešnosti uporabljenih metod in testiranjem metod pri resničnem trgovanju. Na lastnem simulatorju smo preverili, kateri parametri najbolj vplivajo na uspešnost trgovanja. Implementirali smo še nekaj izboljšav, ki so omenjene v sorodni literaturi. Nazadnje smo v poglavju 6.6 preizkusili tudi metode združevanja napovedi samostojnih napovednih modelov in pokazali, da je združevanje napovedi večih modelov bolj uspešno od napovedi posameznih modelov.

Magistrska naloga je razdeljena na sedem poglavij. V poglavju 2 je opisan valutni trg in značilnosti trgovanja na valutnem trgu. V poglavju 3 so na kratko predstavljene uporabljene množice podatkov in njihovo pridobivanje. Poglavje 4 je razdeljeno na dva dela. V prvem so opisane uporabljene metode predstavitve besedil, v drugem delu pa so opisane uporabljene metode strojnega učenja. V poglavju 5 je opisano preverjanje uspešnosti metod predstavitve besedil pri napovedovanju sentimenta s pripadajočimi eksperimentalni rezultati. V poglavju 6 je opisano simulacijsko okolje za posnemanje resničnega trgovanja, nato pa so predstavljeni rezultati obsežnega testiranja implementiranih algoritmov in njihovih parametrov, ter njihova interpretacija. Na koncu poglavja so opisane tudi tehnike združevanja napovedi in njihovi rezultati, ter nekaj najboljših algoritmov, ki so simulacijo končali z do 250% vrednostjo začetnih sredstev. Nazadnje so v poglavju 7 predstavljene sklepne ugotovitve in predlogi za nadaljnje raziskovanje.

Poglavje 2

Valutni trg in trgovanje

Valutni trg najpogosteje poimenujemo kar Forex ali FX, kar v prevodu pomeni mednarodna menjava med valutami (angl. *Foreign Exchange*). Valutni trg je največji ter najlikvidnejši finančni trg na svetu. Na valutnem trgu se opravlja menjava ene valute v drugo, kar imenujemo posel.

Zapis valutnih parov je standardiziran že od leta 1973 in sicer s standardom ISO-4217. Standard določa tričrkovni zapis posamezen valute, valutni par pa je lahko predstavljen z zapisom AAABBB ali AAA/BBB. AAA pravo osnovna valuta, BBB pa obratna valuta. Na spletu¹ je dostopen prosto dostopen seznam večine valut, za potrebe tega magistrskega dela pa smo se omejili le na valutni par EURUSD.

Valutno trgovanje poteka pet delovnih dni na teden. Trgovalni teden se začne v nedeljo ob 20.00, konča pa v petek ob 20.00 po greenwiškem osrednjem času (GMT). Trgovanje ločimo po tem, katere finančne institucije so odprte po svojem lokalnem času. Trgovalni dan se prične ob 20.00 GMT v Wellingtonu, Nova Zelandija, kjer je po lokalnem času ura 08.00. Nadaljuje se s pričetkom delovnega dne v Tokiu ob 22.00 GMT, eno uro pred zaprtjem azijskih finančnih ustanov pa se prične delovni dan v Evropi, pet ur za tem pa še delovni dan v New Yorku in s tem odprtje ameriškega finančnega trga. Predstavljena obdobja so zgolj okvir, po katerem delujejo lokalne fi-

¹<http://www.xe.com/iso4217.php>

nančne institucije, posameznik pa lahko trguje kadarkoli je trg odprt, torej od ponedeljka do petka.

Po ocenah Foreign Exchange Committee [16] je povprečni dnevni promet v aprilu 2014 znašal 811 milijard USD. Najbolj prometni valutni pari so EURUSD (25% dnevnega prometa), USDJPY (16%), GBPUSD (10%) in CADUSD (8%). Največ transakcij se opravi v Veliki Britaniji (37%), ZDA (18%) in na Japonskem (6%). Največje svetovne banke, Deutsche Bank, Citi Bank, Barclays Capital, UBS AG in HSBC, opravijo 53% vseh dnevnih transakcij, ostalo pa prispevajo še centralne banke, ki trgujejo z namenom uravnavanja inflacije in zalog denarja, ter drugi vlagatelji.

Osnovno mersko enoto za količino transakcij izbranega valutnega para predstavlja 100.000 enot osnovne valute in ji pravimo paket (ang. *lot*). V osnovi bi za nakup v višini 100.000 EUR potrebovali 100.000 EUR, vendar banke trgovalcem ponujajo trgovanje s pomočjo finančnega vzvoda. Finančni vzvod (angl. *leverage*) je mehanizem, kjer je trgovalcu (v primeru uporabe finančnega vzvoda v razmerju 1:100) potrebno plačati le $\frac{1}{100}$ vrednosti posla. S tem lahko trgovalec opravlja posle, ki so vredni več kot ima sredstev na računu. Koristi ima tudi posrednik, saj s tem zasluži 100-krat višjo provizijo z razmikom med nakupno in prodajno ceno. Banka potrebuje garancijo, da trgovalec ne bo naredil izgube, ki je ne more pokriti, zato zahteva minimalna rezervirana sredstva (angl. *margin*) v višini vrednosti odprtih poslov. Poleg rezerviranih sredstev mora imeti trgovalec na računu tudi prosta sredstva, ki jih banka potrebuje za poplačilo negativnega stanja poslov. Če negativno stanje odprtih poslov preseže vrednost rezerviranih sredstev, banka nemudoma zapre odprte posle (angl. *margin call*) in jih poplača z rezerviranimi sredstvi.

2.1 Opis trgovanja

Valutno trgovanje se izvaja elektronsko preko različnih trgovalnih platform. Vsak uporabnik trgovalne platforme ima odprt račun pri eni od bank, ki

sodelujejo pri valutnem trgovanju, ter začetna sredstva računa. V trgovalni platformi lahko odpiramo posle, kar pomeni da izberemo valutni par, količino enot in smer, v katero odpremo posel. Navadno trgujemo po trenutni ceni menjalnega tečaja, možno pa je odpreti tudi naročila, ki se izvedejo ob določenih pogojih. Smer posla je lahko nakup ali prodaja. Nakup enot tečaja EURUSD pomeni, da z EUR kupujemo USD, ter obratno, prodaja enot tečaja EURUSD pomeni nakup EUR z USD. Menjalni tečaj je predstavljen z dvema vrednostima, nakupna (angl. *ask*) in prodajna (angl. *bid*) cena. Razlika med nakupno in prodajno ceno se imenuje razmik (angl. *spread*) in je vir zaslužka za banke, ki ponujajo trgovanje. Običajen razmik je odvisen od posamezne banke, običajno je razmik pri tečaju EURUSD enak 0.0003. Vrednosti tečajev se zapisujejo na eno desetisočinko natančno (tj. 0.0001), kar imenujemo točka (angl. *pip*).

Vrednost menjalnega tečaja se določa glede na trenutno ponudbo in povpraševanje. Določajo jo posredniki trgovanja (angl. *brokers*). Celotna struktura trgovalcev se deli na več nivojev, ločijo pa se po kupni moči in posledično velikosti razmika med nakupno in prodajno ceno.

Večina vlagateljev na valutnem trgu je špekulativnih. Po nekaterih ocenah je takih 80%. V to skupino spadajo tvegani skladi, naložbena podjetja in mali vlagatelji. S pojavom spletnih orodij je trgovanje postalo vedno bolj dostopno malim vlagateljem, zato v zadnjem času število njih narašča [9]. Razvoj orodij za elektronsko trgovanje je bankam znižal stroške transakcij, kar se neposredno vidi v zmanjšanju razmika med nakupno in prodajno ceno v zadnjih letih. Banke so v zadnjih nekaj letih s ponudbo naprednih orodij in aplikacijskih programskih vmesnikov (API) omogočile razvoj elektronskih trgovalnih sistemov, ki jih uporablja vedno več trgovalcev.

Finančni analitiki so glede vpliva na gibanje valutnih tečajev na dveh polih. Prvi trdijo, da je gibanje tečajev možno napovedati s pomočjo vzorcev iz preteklosti, kar imenujemo tehnična analiza. Za lažje odkrivanje vzorcev so v preteklosti izdelali mnogo predstavitev, ki jim pravimo tehnični indikatorji. V tem magistrskem delu smo izdelali nekaj napovednih modelov s pomočjo

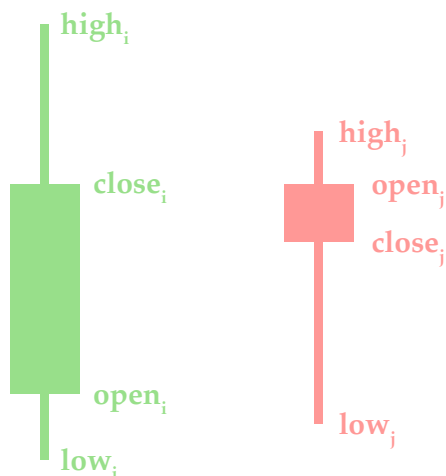
tehničnih indikatorjev, da bi preverili, kako dobro je možno napovedati gibanje tečaja EURUSD. Drugi pol finančnih analitikov trdi, da je gibanje na valutnem trgu odvisno od predpostavk in mnenja družbe o tečaju, kar imenujemo temeljna analiza. Ti analitiki se trudijo zaznati mnenje družbe o tečaju in trgujejo v skladu z njim. Ponavadi trgujejo ob odmevnih družbenih dogodkih, ki vplivajo na predpostavke o tečaju. To so objave različnih državnih agencij (o brezposelnosti, produktivnosti, izvozu, ipd.), aktualni politični dogodki, ter celo večje nesreče ali naravne katastrofe. Predpostavke in mnenje družbe o tečaju je po našem mnenju možno zaznati na družbenih omrežjih, zato smo za vir teh informacij izbrali največje in javno družbeno omrežje Twitter (podobno so izbrali tudi [1, 6, 7]), s katerim smo preizkusili kako dobro lahko s pomočjo temeljne analize napovemo gibanje tečaja EURUSD.

2.2 Tehnični indikatorji

Naloga tehničnih indikatorjev [17] je prikaz vzorcev, trendov in pravil v podatkih o gibanju finančnih tečajev iz preteklosti. Indikatorji so poimenovani v angleščini, od tod tudi kratice, pomen pa je opisan v opisu indikatorja.

Grafi tečajev so sestavljeni iz t.i. japonskih svečnikov (primer svečnika je na sliki 2.1). Vsak svečnik predstavlja časovni interval (minuta, ura, dan, teden), sestavljen pa je iz vrednosti $open_t$, $close_t$, $high_t$ in low_t (vrednost ob začetku časovnega intervala, ob koncu časovnega intervala, najvišja vrednost v časovnem intervalu in najnižja vrednost v časovnem intervalu), kjer t predstavlja čas začetka intervala. Svečnik je pozitiven, če je vrednost ob koncu intervala $close_t$ večja od začetne vrednosti intervala $open_t$ in negativen, če je vrednost tekom izbranega časovnega intervala padla. Vsak interval vsebuje še informacijo o količini trgovanja - $volume_t$, ki ni grafično predstavljena s svečnikom. Te vrednosti so osnova za izračun vrednosti tehničnih indikatorjev in jih uporabljamo v njihovih formulah v nadaljevanju. Spodnji opis indikatorjev smo povzeli po spletni strani².

²<http://ta.mql4.com/indicators>



Slika 2.1: Primer pozitivnega (zelen) in negativnega (rdeč) svečnika, ki prikazujeta gibanje vrednosti tečaja v izbranem intervalu (npr. dnevni ali uri)

AO - Awesome Oscillator

Tehnični indikator AO_t predstavlja povprečno vrednost najvišje in najnižje vrednosti tečaja v intervalu dneva t , ki je definiran z enačbo

$$AO_t = \frac{high_t + low_t}{2} .$$

AC - Acceleration Deceleration

Tehnični indikator AC_t^m , ki v prevodu pomeni pospešek in pojemek, je vrednost, ki je izračunana z drsečim povprečjem tehničnega indikatorja AO_t in je definiran z enačbo

$$AC_t^m = MA_t(AO, m) ,$$

kjer je m širina drsečega okna in t trenuten dan. Funkcija $MA_t(X, m)$ izračuna drseče povprečje vrednosti indikatorja X s širino drsečega okna m .

AD - Accumulation Distribution

Vrednost tehničnega indikatorja AD_t predstavlja količino trgovanja v času t , uteženo z vrednostmi tečaja. Določa jo enačba

$$AD_t = (2close_t - low_t - high_t) \frac{volume_t}{high_t - low_t + AD_{t-1}} ,$$

kjer so $high_t$, low_t , $open_t$ in $close_t$ vrednosti valutnega tečaja (najvišja, najnižja, začetna in končna vrednost) v izbranem časovnem intervalu (v našem primeru smo uporabili en dan).

ATR - Average True Range

Vrednost indikatorja ATR_t^m izračunamo iz drsečega povprečja indikatorja TR_t , ki je določen z

$$TR_t = \max \left(\frac{high_t - low_t}{1}, \frac{close_{t-1} - low_t}{1}, \frac{high_t - open_{t-1}}{1} \right) ,$$

kjer so $high_t$, low_t , $open_t$ in $close_t$ vrednosti tečaja predstavljene na začetku tega odseka. Indikator ATR_t^m je določen z enačbo

$$ATR_t^m = MA_t(TR, m) ,$$

kjer je m širina drsečega okna in $MA_t(TR, m)$ funkcija, ki izračuna drseče povprečje indikatorja TR v obdobju zadnjih m dni. V tem delu uporabljamo 20-dnevno drseče povprečje, oz. ATR_i^{20} .

SMA - Simple Moving Average

Vrednost tehničnega indikatorja SMA_t v dnevu t je povprečje zadnjih m vrednosti absolutne spremembe tečaja, določene z

$$CNG_{i,j} = close_j - open_i .$$

Vrednost tehničnega indikatorja SMA_t^m je določena z enačbo

$$SMA_t^m = MA_t(CNG_{t-1,t}, m) .$$

MOM - Momentum

Tehnični indikator $MOM_{i,j}$ je indikator, ki prikazuje spremembo vrednosti tečaja v času med i in j , kjer $i < j$. Določen je z enačbo

$$MOM_{i,j} = close_j - open_i .$$

RSI - Relative Strength Index

Tehnični indikator RSI_t^m je indikator, ki upošteva zaključne vrednosti posameznih časovnih intervalov (svečnikov). Razmerje med povprečno vrednostjo pozitivnih svečnikov $P_{t,t-m}$ in povprečno vrednostjo negativnih svečnikov $N_{t,t-m}$ v časovnem oknu dolžine m , oz. v intervalu med dnevi $t - m$ in vključno t , je določeno z

$$RS_t^m = \frac{P_{t,t-m}}{N_{t,t-m}} .$$

Vrednost tehničnega indikatorja določa enačba

$$RSI_t^m = 100 - \frac{100}{1 + RS_t^m}$$

2.2.1 Uporaba

S pomočjo tehničnih indikatorjev opišemo dogajanje na valutnem trgu. Definirali smo dve množici atributov, ki jih uporabljamo v nadaljevanju, in smo ju enostavno poimenovali *tehnični1* in *tehnični2*. Množica atributov *tehnični1* vsebuje vrednosti tehničnih indikatorjev $CNG_{t-1,t}$, $CNG_{t-2,t}$, $CNG_{t-3,t}$, $CNG_{t-4,t}$ in $CNG_{t-5,t}$. V množici atributov *tehnični2* so vsi atributi iz *tehnični1*, ter dodatno vrednosti tehničnih indikatorjev AO_t , AC_t^5 , AD_t , ATR_t^{14} , ATR_t^{20} , SMA_t^5 , SMA_t^{14} , SMA_t^{21} , $MOM_{t-14,t}$ in RSI_t^5 . Atributni zapis enega dneva oziroma trgovalnega dneva t je tako vektor vrednosti tehničnih indikatorjev iz posamezne množice atributov.

Poglavje 3

Podatki

V tem poglavju so opisane vse množice podatkov, ki smo jih ustvarili oziroma pridobili za potrebe analize in izvedbo eksperimentov. Za potrebe analize napovedovanja gibanja valutnih tečajev smo zbrali nekaj različnih množic tvitov, ki se razlikujejo po načinu zajema in časovnih intervalih. V nadaljevanju je na kratko opisano, kako smo jih pridobili in za kaj smo jih potrebovali. Na koncu poglavja je opisana tudi množica preteklih vrednosti valutnega para EURUSD, ki smo jo uporabili pri eksperimentih.

3.1 Twitter

Twitter je eno izmed največjih družbenih omrežij. Po podatkih družbe¹ je mesečno aktivnih več kot 300 milijonov uporabnikov Twitterja. Večina objav uporabnikov je javnih, zato veliko raziskovalcev te podatke uporablja za raziskovanja o vplivu tvitov na različnih področjih. Eno najbolj raziskovanih je finančno področje in vpliv tvitov na gibanje različnih finančnih instrumentov [1, 2, 3, 4, 5], zasledili pa smo tudi raziskave v smeri vpliva na prodajo avtomobilov [6], izide volitev [7] in pojav epidemij [8].

Posamezna objava na omrežju Twitter se imenuje tvit (angl. *Tweet*). Tvit je omejen na 140 znakov, večinoma vsebuje veliko pogovornega je-

¹<https://about.twitter.com/company>



Slika 3.1: Primer tvita, ki vsebuje oznako teme, omembo Twitter uporabnika in spletno povezavo.

zika, okrajšav, slengovskih besed, lahko pa vsebuje tudi spletne povezave. Omembe drugih uporabnikov se začnejo z znakom @, enobesedna oznaka teme pa se začne z znakom #. Primer tvita je prikazan na sliki 3.1. Za uporabo tvitov pri napovedovanju je potrebno besedilo v tvitih primerno očistiti in predstaviti z atributi. Uporabljene metode predstavitve besedil so opisane v poglavju 4.1.

Podjetje Twitter je leta 2009 predstavilo funkcionalnost združevanja uporabniških računov Twitter v seznane. Seznam je preprosto množica uporabnikov Twitterja, ki jo je ustvaril nek tretji uporabnik. Seznam uporabniških računov služi boljši preglednosti nad dogajanjem na Twitterju, saj prikazuje le tvite uporabnikov v seznamu. Zelo pogosti so sezname politikov, igralcev, športnikov, zasledimo pa tudi sezname, ki vsebujejo uporabniške račune finančnih ustanov, javnih ustanov, medijev, ipd. Ta funkcionalnost nam je olajšala pridobivanje tvitov, ki so povezani s finančnim dogajanjem. V nadaljevanju opisane množice tvitov smo omejili zgolj na angleški jezik, ki je med tviti najbolj zastopan [18].

3.1.1 Množica tvitov s pripadajočim sentimentom

Delovanje implementiranih metod predstavitev besedil smo najprej ocenili na enostavnem problemu klasifikacije tvitov v pozitiven ali negativen razred, ki predstavlja sentiment tvita (opisano v poglavju 5). Uporabili smo že izdelano

množico podatkov, ki so jo ustvarili Go in sodelavci [19]. Množica vsebuje 1.6 milijona tvitov v angleškem jeziku, označenih s sentimentom tvita. Sentiment tvita je lahko le pozitiven ali negativen (razred 1 ali 0). Oba razreda v množici tvitov sta predstavljena z enakim številom primerov (0.8 milijona tvitov za posamezni razred).

Avtorji so množico tvitov pridobili s pomočjo emotikonov. Omejili so se na izbrano množico emotikonov in jih uvrstili v pozitivne ali negativne emotikone. Emotikoni :), :D in XD so pozitivni, emotikoni :(, :(in :/ pa so negativni. Izbrali so samo tvite, ki so vsebovali označene emotikone. Tvitom s pozitivnimi emotikoni so določili pozitiven razred, tvitom z negativnimi emotikoni pa negativen razred - sentiment. To množico tvitov smo uporabili v prvem delu eksperimentov, ki so opisani v poglavju 5.

3.1.2 Naključni vzorec tvitov

Organizacija Archive team² se ukvarja z ustvarjanjem arhiva spletnih vsebin, med drugim arhivirajo tudi tvite. Vsak mesec objavijo The Twitter Stream Grab, ki vsebuje približno 130 milijonov naključnih tvitov, kar je v približno 50 tvitov na sekundo.

Po nekaterih ocenah to predstavlja zgolj 1% vseh tvitov, vendar točne informacije o številu tvitov podjetje Twitter ne izda. Zbrali smo tvite za mesec marec in april 2015, torej skupno več kot 200 milijonov tvitov. Množico zbranih tvitov smo uporabili pri izdelavi lastnega modela za transformacijo besedil z word2vec vektorji besed, ki je opisan v poglavju 4.1.3.

3.1.3 Tviti izbranih uporabnikov

Cilj magistrskega dela je izdelava napovednih modelov, ki sledijo principom temeljne analize oziroma, z drugimi besedami, za strojno učenje uporabljajo besedilne attribute. Temeljna analiza pravi, da so spremembe na finančnih trgih odraz družbenega mnenja in dogajanja. Za zajem družbenega dogajanja

²<https://archive.org/>

smo izbrali družbeno omrežje Twitter, ki je, prav tako kot valutno trgovanje, globalno.

Za prve napovedne modele smo uporabili naključni vzorec tvitov iz prejšnjega odstavka. Kmalu smo ugotovili, da naključni vzorec vsebuje preveč tvitov, ki nimajo nikakršne povezave s finančnim svetom in predstavljajo šum pri uporabi z našim problemom - napovedovanjem gibanja valutnih tečajev. Iskali smo boljši način zajema tvitov, ki bi po možnosti vsaj bežno spadali v kategorijo *finance*. V ta namen smo uporabili sezname uporabniških računov. Iskali smo tak seznam, ki bi vseboval uporabnike Twitterja iz finančnega sveta.

Kot najbolj primeren in dovolj velik uporabniški seznam smo izbrali Money Finance Invest³. Avtor, Twitter uporabnik *howtoenjoyldn*, je v seznam zbral 640 uporabniških računov Twitter, bančnike, investitorje, finančne novinarje ali uporabnike, kako drugače povezane z globalnimi financami.

Twitter v splošnih pogojih dovoljuje zbirati le zadnjih 3200 tvitov posameznega uporabnika. To nam ni predstavljalo posebno velike težave, saj je omejitev dovolj visoka, da smo za večino uporabnikov pridobili vse njihove tvite iz preteklih petih let. Skupaj smo za obdobje petih let zbrali 1,224,679 tvitov, ki jih je objavilo 640 Twitter uporabnikov. To v povprečju predstavlja 28.3 tvitov na dan.

3.2 Pretekle vrednosti valutnih tečajev

Podatke o vrednostih valutnih tečajev smo pridobili na spletni strani švicarske banke Dukascopy⁴, ki je ena od bank ponudnic Forex trgovanja v Evropi. V tem magistrskem delu smo se omejili na tečaj EURUSD. Pridobili smo dnevne vrednosti tečaja EURUSD, in sicer za enourni interval in dnevni interval med 1.11.2010 in 30.6.2015.

Podatki so sestavljeni iz zapisov vrednosti začetne, najvišje, najnižje in

³<https://twitter.com/howtoenjoyldn/lists/money-finance-invest>

⁴<https://www.dukascopy.com/free/candelabrum/>

končne vrednosti tečaja v izbranem intervalu. Te štiri vrednosti označimo z $open_i$, $high_i$, low_i in $close_i$, in so osnova za prikaz grafa vrednosti valutnega tečaja in tudi osnova za izračun vrednosti tehničnih indikatorjev, ki so opisani v poglavju 2.2. Nekaj primerov zapisov vrednosti valutnega tečaja:

Time	Open	High	Low	Close	Volume
21.04.2014 22:00:00.000	1.37945	1.37956	1.37907	1.37910	3076.89
21.04.2014 23:00:00.000	1.37913	1.37947	1.37912	1.37916	1048.24
22.04.2014 00:00:00.000	1.37916	1.37943	1.37899	1.37907	1742.02
22.04.2014 01:00:00.000	1.37909	1.37950	1.37900	1.37933	2143.72

Gibanje tečaja EURUSD v obdobju med 1.7.2010 in 1.7.2015 je prikazano na sliki 3.2. Povprečna dnevna sprememba tečaja v tem obdobju je 0.00519 USD oz. 51.9 točke. Največja sprememba tečaja v enem dnevu je bila 339.5 točke, največja negativna dnevna sprememba pa 307.7 točke. 49.8% dni v izbranem obdobju se konča s pozitivno spremembo tečaja, 51.2% dni pa z negativno spremembo tečaja.



Slika 3.2: Gibanje tečaja EURUSD v obdobju med 1.7.2010 in 1.7.2015

Poglavje 4

Uporabljene metode

V tem poglavju so predstavljene uporabljene in implementirane metode za predstavitev besedil, v našem primeru tвитov. Da lahko besedila uporabimo kot učne podatke pri metodah strojnega učenja, jih moramo predstaviti v vektorski oziroma atributni obliki. V poglavju 4.1 predstavimo dve znani in pogosto uporabljeni predstavitvi besedil - vrečo besed in vrečo nizov. Dodatno smo raziskali predstavitev besedil z word2vec vektorji besed, ki so opisani v poglavju 4.1.3. V nadaljevanju poglavja so opisane uporabljene metode strojnega učenja (poglavje 4.2.1) in tehnika sprotnega izbora atributov (poglavje 4.2.2), ki smo jo zasledili v sorodni literaturi. Nazadnje sta opisani še dve uporabljeni tehniki združevanja napovedi.

4.1 Predstavitev besedilnih dokumentov

Pri strojnem učenju primere predstavimo z vektorji atributov. Če so primeri besedilni dokumenti ali množice tvitov, jih prav tako predstavimo v atributnem zapisu, ki je primeren za metode strojnega učenja, s katerimi lahko gradimo napovedne modele. Za pretvorbo v atributni zapis uporabimo eno od metod predstavitev besedilnih dokumentov, ki so opisane v nadaljevanju tega poglavja.

Še pred izbiro ustrezne predstavitve besedil je potrebno odstraniti nepolnomenne besede (angl. *stopwords*). V našem primeru so besedilni dokumenti tviti. Zanje veljajo posebne značilnosti. Spisali so jih uporabniki Twitterja, lahko vsebujejo spletne povezave, omembe drugih uporabnikov, ter eno ali več oznak tem s ključnimi besedami (angl. *hashtags*). Po naši oceni so spletne povezave in uporabniška imena zgolj šum, zato jih prav tako, kot nepolnomenne besede, odstranimo.

Tviti so v večini napisani v pogovornem jeziku, vsebujejo veliko okrajšav, slenga in pogosto tudi emotikone. Te značilnosti otežijo predstavitev besedila, zato predlagamo tudi uporabo transformacije z word2vec vektorji besed [20], ki v določenih primerih sopomenke ali nepravilno zapisane besed z istim pomenom predstavi s podobnim vektorjem. Delovanje in uporaba word2vec vektorjev besed je podrobneje opisano v poglavju 4.1.3.

4.1.1 Vreča besed

Vreča besed (angl. *bag of words*) je ime za predstavitev besedila, kjer vsak dokument d_i predstavimo z vektorjem

$$d_i = (f_{i,1}, f_{i,2}, \dots, f_{i,n}) ,$$

ki je sestavljen iz frekvenc $f_{i,j}$ besede w_j v dokumentu d_i . Primer transformacije enega dokumenta z vrečo besed je prikazan na sliki 4.1. Če imamo na primer 10 dokumentov in skupaj vsebujejo 1000 unikatnih besed, bomo s to metodo besedila predstavili z matriko atributov dimenzij (10, 1000).

John likes to watch movies. Mary likes movies too.

john	likes	watch	movies	mary	thing
1	2	1	2	1	0

Slika 4.1: Primer predstavitve besedila z vrečo besed

Izboljšava predstavitve besedil z vrečo besed je preslikava TF-IDF (angl. *term frequency, inverse document frequency*), ki posamezno besedo v dokumentu oziroma atribut uteži glede na frekvenco pojavljanja besede v dokumentu in frekvenco pojavljanj besede v celotnem naboru dokumentov. Preslikava TF-IDF besede t v dokumentu d pri naboru vseh dokumentov D je določena kot

$$\text{tfidf}(t, d, D) = \text{tf}(t, d) \cdot \text{idf}(t, D) = \text{tf}(t, d) \cdot \log \frac{N}{|\{d \in D : t \in d\}|},$$

kjer je $\text{tf}(t, d)$ število pojavitev besede t v dokumentu d , vrednost $\text{idf}(t, D)$ pa je logaritmično utežena vrednost količnika med številom vseh dokumentov N in številom tistih dokumentov, ki vsebujejo besedo t , kar je definirano z $|\{x \in D : t \in x\}|$.

Uporaba preslikave TF-IDF zmanjša težo besed, ki se uporabljajo v večini dokumentov in poveča težo besedam, ki se pojavijo zgolj v nekaj dokumentih. Že med implementacijo predstavitve besedil z vrečo besed smo na podlagi testnih poskusov ugotovili, da uporaba preslikave TF-IDF tipično izboljša klasifikacijsko točnost. Zaradi velikega števila drugih nastavitev in parametrov metod smo se odločili, da v vseh eksperimentih, ki za predstavitev besedila uporabljajo vrečo besed, v nadaljevanju uporabimo tudi preslikavo TF-IDF.

Implementacija vreče besed v pythonovski knjižnici scikit-learn omogoča tudi uporabo besednih zvez. Tako lahko besedno zvezo iz dveh ali več besed predstavimo z enim atributom. Če želimo namesto posamezne besede

upoštevati besedno zvezo, lahko v knjižnici scikit-learn uporabimo parameter $ngrams = (1, 2)$, ki določa, katere dolžine besednih zvez bomo uporabili. Zapis $(1, 2)$ pomeni, da so atributi predstavljeni iz ene samostojne besede ali dveh besed. S testnimi poskusi na problemu napovedovanja sentimenta smo ugotovili, da najvišjo klasifikacijsko točnost doseže vreča besed, ki ima za attribute samostojne besede, ter besedne zveze iz dveh in treh besed, zato smo v knjižnici scikit-learn uporabili parameter $ngrams = (1, 3)$. V praksi se uporablja tudi omejevanje uporabe besed, ki imajo frekvenco manjšo od določene vrednosti. V eksperimentih smo upoštevali le besede, ki se v besedilu pojavijo vsaj šestkrat.

4.1.2 Vreča nizov

Vreča nizov je metoda za predstavitev besedil, podobna vreči besed. Edina razlika je v tem, da osnovna enota, ki tvori atribut, ni beseda, ampak podniz znakov. Podniz znakov tipično pridobimo tako, da se z drsečim oknom določene dolžine pomikamo čez besedilo. Pri vreči nizov prav tako lahko nastavimo minimalno in maksimalno željeno dolžino nizov s parametrom scikit-learn knjižnjice $ngrams$. Primer transformacije kratkega besedila z uporabo parametra $ngrams = (2, 4)$ je prikazan na sliki 4.2. S testnimi poskusi na problemu napovedovanja sentimenta smo ugotovili, da je smiselno uporabiti nize dolžine $(2, 6)$. V vseh eksperimentih smo v povezavi z vrečo nizov uporabili parameter $ngrams = (2, 6)$ in nad dobljeno predstavitvijo potem uporabili preslikavo TFIDF. Prav tako smo pri uporabi te predstavitve besedila upoštevali le nize, ki se v besedilu pojavijo vsaj šestkrat.

4.1.3 Predstavitev z word2vec vektorji besed

Značilnosti tvitov, kot so razširjena uporaba slenga, okrajšav in pravopisne napake, otežijo pravilno predstavitev z vrečo besed. Vreča nizov lahko določene primere predstavi z ujemačimi se podnizi, vendar se tudi pri njej pojavi problem: ne moremo združiti podobnih besed z istim pomenom v en

Nor agai is there anyone who loves or pursues or desires the ...

agai	gain	dust	no	or	the	per	taly
1	1	0	1	3	2	0	0

Slika 4.2: Primer predstavitve besedila z vrečo nizov dolžine dveh, treh ali štirih znakov.

atribut. Na primer, besede *grandma*, *granma*, *grandmother* pomenijo isto, zapisane pa so drugače.

Dodaten razlog za iskanje alternativnih predstavitev besedil je tudi visoka prostorska zahtevnost metod vreče besed in vreče nizov. Zato predlagamo alternativno predstavitev besedil, ki dokument prav tako predstavi v vektorski obliki in sicer z uporabo word2vec vektorjev besed.

Pristop word2vec

Mikolov in sodelavci [20] opišejo postopek, pri katerem iz velike množice besedil pridobijo vektorje besed, ki ohranjajo medsebojne podobnosti in povezave. Za učenje vektorjev besed uporabljajo umetne nevronske mreže (angl. *artificial neural networks* - ANN).

Nevronska mreža je metoda strojnega učenja, ki posnema delovanje človeških oziroma živalskih možganov. Sestavljena je iz nevronov, razporejenih v sloje. Najenostavnejše nevronske mreže so zgrajene iz enega sloja, bolj komplicirane pa jih lahko vsebujejo več. Vsaka nevronska mreža ima vhodni sloj in izhodni sloj, dodatne sloje (med vhodnim in izhodnim) pa imenujemo skriti sloji. Nevroni različnih slojev so med seboj povezani s povezavami, ki posnemajo živčne povezave v možganih, po katerih se pošiljajo živčni signali. Neuron je aktiven in po izhodni povezavi pošilja signal, ko je vsota njegovih vhodnih signalov dovolj velika.

Parametri nevronske mreže so uteži vhodov na posameznih nevronih in

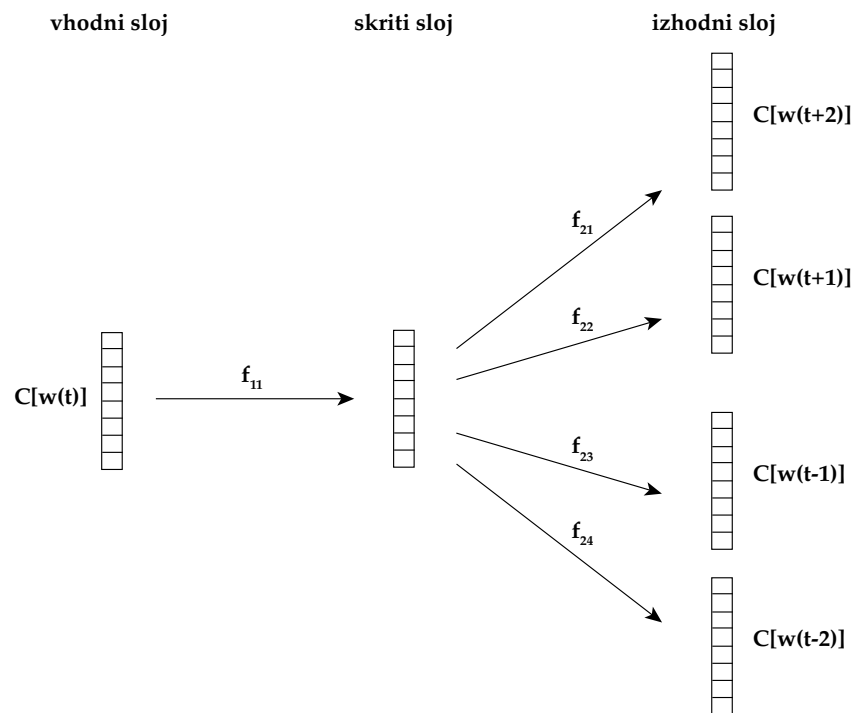
pragovi, pri katerih se posamezni nevroni aktivirajo. Ti parametri se med učenjem na učnih podatkih spreminjajo, dokler nevronska mreža ne zna sama dovolj dobro rešiti problema, na katerem se uči.

Nevronska mreža, ki jo predlagajo Mikolov in sodelavci v [20] je enostavna nevronska mreža z enim skritim slojem (slika 4.3). Poleg nevronske mreže avtorji definirajo tudi matriko C , ki jo spreminjamo tekom učenja nevronske mreže. Matrika C je sestavljena iz $|V|$ vektorjev besed C_i , ki so dimenzije n . Slovar besed V , ki jih upoštevamo, mora biti določen že pred učenjem nevronske mreže in ga algoritem pridobi z enkratno iteracijo skozi učno besedilo.

Nevronska mreža na sliki 4.3 je torej sestavljena tako, da za vektor besede w_t napove najbližje štiri sosednje vektorje besed $C[w_{t-2}]$, $C[w_{t-1}]$, $C[w_{t+1}]$ in $C[w_{t+2}]$. Za napovedovanje je potrebno učenje nevronske mreže na množici besedil. Učenje se začne z naključnim izborom vrednosti v matriki C in ostalih parametrov nevronske mreže. Nato se učenje nevronske mreže izvaja v korakih, kjer se za vsak korak z drsečim oknom premikamo čez besedilo. V vsakem koraku na vhod nevronske mreže postavimo vektor besede w_t , da naredi napovedi s trenutnimi parametri oziroma na izhodnih nevronih določi vektorje sosednjih besed. Nato se na podlagi resničnih vektorjev sosednjih besed izračunajo napake, s katerimi nevronska mreža popravi svoje parametre oziroma transformacijske funkcije $f_{i,j}$ in posodobi matriko vektorjev besed C . Temu postopku učenja pravimo vzvratno razširjanje napake, ali angl. *backpropagation*.

Rezultat učenja nevronske mreže je torej matrika vektorjev besed C , ki v n dimenzijah ohranja netrivialne podobnosti med besedami, ki se pojavijo v tekstu. Skupaj v vektorskem prostoru se pojavljajo vektorji sopomenk, vektorji imen držav in prav tako vektorji besed z istim korenem in različnimi končnicami. Blizu v vektorskem prostoru vektorjev word2vec so tudi vektorji nepravilno zapisanih besed.

Avtorji pristopa word2vec kot veliko prednost njihove metode izpostavljajo enostavne numerične operacije nad vektorji besed. Najbolj znan pri-



Slika 4.3: Arhitektura nevronske mreže za učenje word2vec vektorjev besed (matrike C)

mer [21] je vsota vektorjev besed “king” - “man” + “woman”, ki predstavljajo vektor, kateremu najbližja beseda je “queen”. Podobno v članku [22] izpostavijo izračun “madrid” - “spain” + “france” = “paris”. V primeru dovolj velike večjezične množice učnih besedil se vektorji iste besede v različnih jezikih ravno tako pojavljajo skupaj [22].

Učenje transformacijskega modela word2vec potrebuje veliko množico besedil in je nenadzorovano. Tipično se za učenje uporabi besedila iz podobnega konteksta ali celo besedila, ki jih uporabljamo v problemu, ki ga rešujemo.

Uporaba preslikave word2vec

Vektorje besed word2vec smo uporabili za transformacijo besedil v atributne opise. Besede v tvitih smo predstavili z word2vec vektorji. Besede, ki jih transformiramo v vektorje, morajo biti prisotne v slovarju V , kar pomeni, da ne moremo uporabiti vseh besed v tvitih. Za naše potrebe je transformacija v word2vec vektorje besed uporabna zaradi enostavne in hitre preslikave besed v nizko-dimenzionalen prostor, ter ohranjanja povezav med besedami. Posamezno besedo torej predstavimo z vektorjem fiksne dolžine (uporabljen parameter $n = 300$).

Posamezen primer, ki ga moramo predstaviti v vektorski obliki, je v našem primeru tvit, ki je sestavljen iz več besed, zato je potrebno vektorje besed združiti v neko smiselno predstavitev, tako da en vektor predstavlja celoten tvit. Po zamisli spletnega vodiča uporabe word2vec vektorjev besed¹ smo implementirali dva načina združevanja vektorjev besed v vektor dokumenta. To sta pristopa *povprečje word2vec* in *skupine word2vec*.

povprečje word2vec kot že ime pove, posamezne vektorje besed združi v vektor dokumenta tako, da enostavno izračuna povprečje komponent vektorjev besed.

skupine word2vec pristop združevanja vektorjev besed v vektor dokumenta najprej razvrsti vektorje besed z metodo voditeljev v 100 skupin. Do-

¹<https://www.kaggle.com/c/word2vec-nlp-tutorial/>

kument oziroma množico besed predstavimo z vektorjem frekvenc pojavitev besed posamezne skupine v dokumentu. Vektor dokumenta dobimo tako, da vsaka beseda v dokumentu glasuje oziroma poveča števec skupine besede.

Pri eksperimentih smo uporabili dve različni matriki vektorjev besed. Prvo so izdelali Mikolov in sodelavci [20]. Izdelali so jo z učenjem na množici novic Google News s skupno 100 milijard besed, matrika pa vsebuje 1.4 milijona vektorjev besed. Matrika word2vec vektorjev besed vsebuje zgolj 60% besed iz tvitov, ki jih uporabljamo za izdelavo napovedi sentimenta tvita (poglavje 5).

Iz dvomesečne množice zajetih tvitov (poglavje 3.1.2) smo naučili novo nevronske mreže z novo word2vec matriko C . Po nekaj poskusih napovedovanja sentimenta tvita smo ugotovili, da je pri uporabi lastno izdelane matrike word2vec vektorjev zastopanost besed v slovarju večja. Besed izven slovarja je bilo le še 27%. Dvomesečna množica zajetih tvitov je tako vsebovala 77% besed, ki so se pojavile v množici tvitov označenih s sentimentom. Mikolov in sodelavci [22] poročajo, da se, v primeru zmanjšanja množice učnega besedila, zmanjša točnost napovedovanja na testnem problemu, ki ga uporabljajo za ocenjevanje kakovosti izdelanih vektorjev word2vec. V vseh eksperimentih zato uporabljamo slednjo, lastno izdelano word2vec matriko vektorjev besed.

Uporabili smo implementacijo word2vec vektorjev besed v pythonovski knjižnici Gensim². Knjižnica omogoča učenje word2vec vektorjev, ki je opisano v prejšnjem razdelku, in enostavno shranjevanje in uporabo matrike C .

²<https://radimrehurek.com/gensim/>

4.2 Napovedovanje

Uporabljene metode strojnega učenja so na kratko opisane v poglavju 4.2.1. V literaturi smo zasledili zanimiv pristop, ki za vsak primer za katerega napovedujemo vrednost razreda posebej izbere le attribute besed, prisotnih v besedilu primera. Ta pristop je opisan v poglavju 4.2.2. Uporabili smo tudi tehnike združevanja napovedi različnih modelov, ki so podrobneje opisane v poglavju 4.2.3.

4.2.1 Metode strojnega učenja

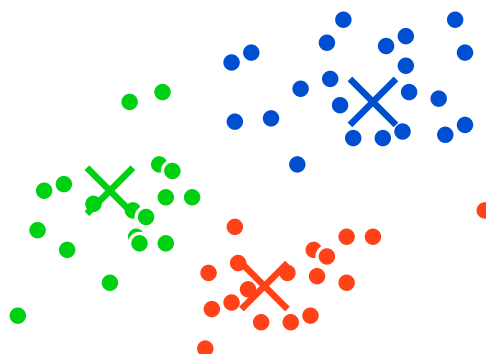
Strojno učenje smo uporabili za napovedovanje gibanja tečajev. V splošnem pristope strojnega učenja delimo na dve večji skupini, nadzorovane in nenadzorovane. Pri nadzorovanem učenju (angl. *supervised learning*) se algoritem uči pravil oziroma izdelava model iz množice učnih podatkov in znanih vrednosti razreda, nato pa na podlagi teh pravil novim primerom napove vrednost razreda. Nenadzorovano učenje (angl. *unsupervised learning*) ne potrebuje znanih vrednosti razreda in učne primere razvrsti v skupine na podlagi njihovih značilnosti in uporabljene metodologije.

Pri napovedovanju prihodnjega gibanja valutnih parov smo uporabili metode strojnega učenja iz obeh skupin (nadzorovano in nenadzorovano učenje). Vse uporabljene metode nadzorovanega strojnega učenja so napovedovale diskretne razrede, zato smo uporabili klasifikacijske različice metod strojnega učenja implementirane v knjižnici scikit-learn³.

Metoda voditeljev

Metoda voditeljev (angl. *K-means*) spada v skupino nenadzorovanih metod strojnega učenja. Uporabnik določi število skupin, na katero želi razdeliti množico primerov. Metoda iz podatkov poišče k skupin in njihovih voditeljev, oziroma središč skupin. Iskanje optimalnega razbitja primerov v skupine se začne z naključno izbranimi voditelji, nato pa algoritem uvrsti vsakega od

³<http://scikit-learn.org/stable/>



Slika 4.4: Razbitje podatkov v tri skupine z metodo voditeljev. S krogi so predstavljeni primeri, križi predstavljajo središča skupin.

primerov v skupino najbližjega voditelja. Po zaključenem uvrščanju primerov se na novo izračunajo središča skupin oziroma voditelji, celoten postopek pa se večkrat ponovi, dokler se lega voditeljev ne spreminja več. Primer rezultata takega razvrščanja z metodo voditeljev s tremi voditelji je prikazan na sliki 4.4.

Najpomembnejši parameter te metode je število skupin. Dodatni parametri so še izbira ustavitvenega kriterija, način izbora začetne pozicije voditeljev in način izračuna razdalj med primeri in središči. Implementacija metode v knjižnici `scikit-learn`, ki jo uporabljamo, privzeto uporablja naprednejšo izbiro začetnih voditeljev, zato tega parametra nismo posebej spreminjali. Metodo voditeljev uporabljamo pri predstavitvi tвитov z `word2vec` vektorji besed, ki je opisana v poglavju 4.1.3.

Logistična regresija

Logistična regresija (angl. *Logistic regression*) je linearen model in spada v družino posplošenih linearnih modelov [23]. Podobna je metodi linearne regresije, vendar nasprotno od pomena imena napoveduje diskretno vrednost razreda. Logistična regresija primerom $x = [x_1, x_2, \dots, x_n]^T$ napove razred s

pomočjo funkcije

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}} ,$$

kjer funkcijo

$$g(z) = \frac{1}{1 + e^{-z}}$$

imenujemo logistična funkcija ali sigmoida. Vloga funkcije $g(z)$ je omejevanje izhoda funkcije $h_{\theta}(x)$ na interval $[0, 1]$.

Model logistične regresije je torej določen s parametri $\theta = [\theta_1, \theta_2, \dots, \theta_n]^T$. Metoda primeru x določi razred tako, da izračuna vrednost funkcije $h_{\theta}(x)$. Vrednost te funkcije je verjetnost, da primer spada v pozitiven razred. Za določitev končnega razreda določimo še prag verjetnosti za pozitiven razred, tipično kar $P > 0.5$. Primeru x določimo pozitiven razred če $h_{\theta}(x) > 0.5$ in negativen razred če $h_{\theta}(x) < 0.5$.

Učenje modela logistične regresije je iskanje takih parametrov θ , da funkcija $h_{\theta}(x)$ čimbolje loči razreda učnih primerov. Postopek učenja izvedemo z gradientno metodo (angl. *gradient descent*), ki minimizira funkcijo napake (angl. *cost function*). To je funkcija, ki meri točnost napovedi trenutnega modela oziroma parametrov θ . Gradientna metoda na začetku učenja naključno izbere parametre θ_i in z iteriranjem skozi učne primere postopoma in z manjšimi koraki popravlja parametre v smeri zmanjševanja napake, ki je določena z odvodom funkcije napake.

Rezultat učenja je optimalen nabor parametrov θ , ki določajo mejo med primeri obeh razredov. Če so učni podatki predstavljeni z velikim številom atributov, se metoda logistične regresije prekomerno prilagodi učnim podatkom. Takrat pravimo, da je model prenasičen (angl. *overfitted*) in točnost napovedovanja novih primerov občutno pade. Prenasičenju modela logistične regresije se lahko izognemo z uporabo regularizacije [24], ki med postopkom učenja zmanjšuje vrednost oziroma težo parametrov θ_i . Preizkusili smo tri različne tipe regularizacije, *l1*, *l2* in *elasticnet*, ki so implementirani v knjižnici `scikit-learn`⁴. Dodatno je potrebno poiskati tudi primeren parame-

⁴http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier.html

ter α , ki določa faktor zmanjševanja θ_i . V knjižnici scikit-learn se pri uporabi logistične regresije privzeto uporablja tip regularizacije $l2$ in parameter $\alpha = 1.0$.

Če želimo z logistično regresijo napovedovati več kot dva diskretna razreda, se moramo poslužiti tehnik kombiniranja večih modelov. Implementacija logistične regresije v knjižnici scikit-learn uporablja tehniko eden proti ostalim (angl. *one vs. rest*). Ta tehnika izdelava r klasifikatorjev, kjer je vsak klasifikator naučen ločevanja primerov med izbranim razredom in preostalih $r - 1$ razredi. Vsak od r klasifikatorjev za posamezen primer napove verjetnost izbora posameznega razreda. Za končno napoved skupine r klasifikatorjev je izbran razred z največjo verjetnostjo. V primeru neodločenega izida se pogosto uporablja naključen izbor izmed razredov z največjo verjetnostjo, lahko pa tudi druge metode [25].

Optimalne vrednosti parametrov regularizacije določimo s testiranjem na problemu napovedovanja sentimenta (poglavje 5.2.1). Logistično regresijo nato uporabimo pri napovedovanju gibanja valutnega tečaja, kjer primere klasificiramo v enega od diskretnih razredov. Metodo uporabimo tudi pri združevanju napovedi večih klasifikatorjev.

Naivni Bayesov klasifikator

Naivni Bayesov klasifikator združuje dva pristopa, ki ju lahko prepoznamo že iz imena. Naivni je zaradi naivne predpostavke pogojne neodvisnosti atributov, Bayesov pa zaradi uporabe Bayesovega izreka.

Predpostavka pogojne neodvisnosti atributov predpostavlja, da so vrednosti atributov a_i pri izbranem razredu c_j med seboj neodvisne in jih lahko zapišemo kot produkt posameznih pogojnih verjetnosti

$$P(a_1, a_2, \dots, a_i | c_j) = \prod_i P(a_i | c_j) ,$$

kjer so a_i vrednosti atributov in c_j vrednosti razreda.

Osnovni Bayesov izrek za pogojno verjetnost $P(A|B)$ je

$$P(A|B) = \frac{P(A)P(B|A)}{P(B)} .$$

Z njim lahko zapišemo tudi pogojno verjetnost razreda c_j pri danih atributih a_i kot

$$P(c_j|a_1, a_2, \dots, a_i) = \frac{P(c_j)P(a_1, a_2, \dots, a_i|c_j)}{P(a_1, a_2, \dots, a_i)}.$$

Po predpostavki pogojne neodvisnosti posameznih atributov a_i in s pomočjo Bayesovega izreka za verjetnost razreda, lahko verjetnost razreda zapišemo z enačbo

$$P(c_j|a_1, a_2, \dots, a_n) = P(c_j) \prod_i^n \frac{P(a_i|c_j)}{P(a_i)},$$

ki jo uporablja naivni Bayesov klasifikator.

Člen $P(a_i)$ je v enačbi konstanta, zato ga lahko izpustimo. Naivni Bayesov klasifikator pri napovedi razreda za podane attribute a_i izračuna verjetnosti vseh možnih razredov c_j , izbrani in napovedani razred

$$\hat{y} = \underset{c_j}{\operatorname{argmax}} P(c_j) \prod_i^n P(a_i|c_j)$$

pa je tisti z največjo verjetnostjo.

Učenje klasifikatorja je aproksimacija verjetnosti $P(c_j)$ na učnih primerih, ki jim pravimo tudi apriorne verjetnosti. Med učenjem lahko izpustimo take attribute, ki za dani razred c_j nimajo vrednosti. Za učenje potrebujemo diskretne vrednosti atributov, zato v primeru zveznih vrednosti poskrbimo za diskretizacijo.

Uporabljali smo klasifikator implementiran v knjižnici scikit-learn in sicer kot meta-učni algoritem pri metodi stacking, ki je opisana v nadaljevanju poglavja.

Metoda naključnih dreves

Metoda naključnih dreves (angl. *random forest*) spada v ansambelske metode. Ansambelske metode izdelajo večje število enostavnih napovednih modelov in njihove napovedi združijo v eno napoved. Točnost ansambelskih metod je ponavadi boljša od točnosti najboljšega posameznega enostavnega modela, pod pogojem, da so enostavni napovedni modeli med seboj različni.

Osnova za metodo naključnih dreves je odločitveno drevo. Vsako odločitveno drevo je sestavljeno iz vozlišč, ki razdelijo primere na dve poddrevesi, glede na izbran atribut in pogoj. Postopek gradnje odločitvenega drevesa oziroma učenje poteka tako, da v vsakem vozlišču drevesa izberemo najbolj informativen atribut učne množice, ki deli učno množico na poddrevesa. Deljenje ponavljamo, dokler v listih dreves ne ostanejo le primeri z istim razredom. Odločanje o razredu novega primera je preprost sprehod skozi strukturo drevesa. Začnemo v izhodišču drevesa in na vsakem vozlišču sledimo pogojem, ki ustrezajo primeru, dokler ne pridemo do listov in tako določimo razred primera.

Različnost enostavnejših modelov dosežemo z vpeljavo naključnega izbora atributov, kot je predlagal Breiman v članku [26]. V postopku izbire atributov ne izbiramo več najbolj informativnega izmed vseh a atributov, ampak izmed naključno izbranih m atributov, kjer je $m < a$.

Pri uporabi metode naključnih gozdov je potrebno izbrati število odločitvenih dreves n in število naključno izbranih atributov m pri gradnji posameznega odločitvenega drevesa. Za naše eksperimente smo izbrali $n = 100$, parameter m pa je v knjižnici scikit-learn privzeto izbran tako, kot v članku [26] priporoča Breiman, in sicer $m = \sqrt{a}$, kjer je a število vseh atributov. Metodo naključnih dreves uporabimo združevanju napovedi samostojnih klasifikatorjev.

4.2.2 Sproten izbor atributov

Množico učnih podatkov sestavljeno iz besedil pretvorimo v vektorsko obliko z eno od transformacij iz poglavja 4.1. Če transformacijo naredimo s pomočjo vreče besed (kot opisano v poglavju 4.1.1), je dobljena matrika redka matrika z a atributi, in le m atributov, $m < a$, ima vrednost različno od 0.

Nassirtoussi in sodelavci v članku [1] opisujejo zanimiv pristop sprotno izdelave napovednega modela, ki pri vsaki posamezni napovedi izbere samo prisotne oziroma neničelne attribute v primeru x . Z izbranimi atributi nato za vsako napoved na novo izdelajo napovedni model z eno od metod strojnega

učenja.

Tak način napovedovanja je časovno bolj zahteven, saj je potrebno za napoved enega primera vsakokrat na novo naučiti napovedni model, vendar Nassirtoussi in sodelavci [1] navajajo, da so tako dosegli večjo klasifikacijsko točnost. Sproten izbor atributov smo implementirali in preizkusili pri simulaciji trgovanja. Eksperimentalni rezultati so podani v poglavju 6.

4.2.3 Združevanje napovedi več modelov

Združevanje napovedi več modelov [27] je v zadnjem času pogosto uporabljena tehnika za izboljšanje klasifikacije. Za razliko od ansambelskih metod, za združevanje napovedi ne potrebujemo vedno istega tipa klasifikacijskega modela, kot na primeru metode naključnih gozdov. V magistrskem delu smo preizkusili dva načina združevanja napovedi, ki sta opisana tudi v [27], enostavno glasovanje modelov in metoda stacking, oba opisana v nadaljevanju poglavja.

Glasovanje

Napoved razreda novemu primeru z uporabo glasovanja poteka tako, da vsi modeli v množici napovedo svoj razred. Končni razred primera je določen z večinsko izglasovanim razredom. V primeru neodločenega izida naključno izberemo enega izmed razredov z največ glasovi.

Metoda stacking

Metoda stacking za združevanje napovedi ločenih klasifikacijskih modelov uporablja dodaten klasifikacijski model, zato za njeno delovanje potrebujemo učno množico primerov, da metodo naučimo pravilnega združevanja samostojnih učnih primerov. Meta-učni algoritem je eden od algoritmov strojnega učenja, tako pa mu pravimo zgolj zato, ker se uči pravilne izbire učnih modelov.

Metoda je sestavljena iz klasifikacijskih modelov, ki vsak zase izdela napoved v obliki verjetnosti pozitivnega razreda. Napovedi verjetnosti klasifikatorjev združimo v vektor, ki predstavlja primer za meta-učenje z izbrano metodo strojnega učenja. Metoda stacking se torej na učnih primerih (napovedanih verjetnostih posameznih modelov) meta-nauči klasificirati primere v diskretne razrede.

Učenje pravil združevanja torej poteka tako, da na prvi polovici primerov vsak od klasifikatorjev izdela napoved verjetnosti pozitivnega razreda, meta-klasifikator pa se na podlagi verjetnosti in resničnih vrednosti razredov nauči združevanja. Klasifikacija novih primerov poteka podobno. Vsi klasifikacijski modeli najprej napovedo verjetnost pozitivnega razreda, nato pa naučen meta-model iz vektorja verjetnosti klasificira primer v končni razred. Eksperimentalni rezultati uporabe metode stacking in glasovanja so predstavljeni v poglavju 6.6.

Poglavje 5

Eksperimentalno vrednotenje uspešnosti napovednih tehnik

Med implementacijo metod za predstavitev besedil smo želeli oceniti njihovo delovanje na preprostem problemu. Napovedovanje gibanja tečajev na podlagi besedil se nam je zdel preveč kompleksen problem za enostavno oceno pravilnosti delovanja omenjenih metod, zato smo v ta namen uporabili problem napovedovanja sentimenta tvitov. Uporabili smo množico podatkov z 1.6 milijoni primerov, tj. 1.6 milijona tvitov [19]. Razred vsakega primera je njegov sentiment (0 - negativen, 1 - pozitiven). Množica podatkov je podrobneje opisana v poglavju 3.1.1.

Želeli smo primerjati klasifikacijsko točnost pri uporabi različnih predstavitev besedil. Preizkušali smo tudi uporabo različnih metod strojnega učenja, vendar smo po nekaj začetnih eksperimentih ugotovili, da med njimi ni posebno velikih razlik v klasifikacijski točnosti. V nadaljevanju tega poglavja so predstavljeni zgolj eksperimenti z uporabo logistične regresije kot metode strojnega učenja in iskanje optimalnih parametrov regularizacije logistične regresije.

5.1 Klasifikacijska točnost

Ker imamo pri napovedovanju sentimenta dvorazreden problem, smo za merjenje uspešnosti uporabili najenostavnejše merilo uspešnosti, ki izhaja iz matrike zmot (angl. *confusion matrix*, prikazana na sliki 5.1), klasifikacijsko točnost. Pozitiven sentiment označimo s pozitivnim razredom, negativen sentiment pa z negativnim razredom. S P označimo število pozitivnih tvitov, z N pa število negativnih tvitov. Število vseh primerov označimo z n in velja $n = P + N$. S TP označimo število pravilno klasificiranih pozitivnih primerov, s FP število napačno klasificiranih pozitivnih primerov, s TN število pravilno klasificiranih negativnih primerov, s FN pa število napačno klasificiranih negativnih primerov.

		Napovedana vrednost	
		pozitivna	negativna
Pravilna vr.	poz.	TP	FN
	neg.	FP	TN

Tabela 5.1: Matrika zmot

Klasifikacijsko točnost označimo z ACC , definiramo pa jo z

$$ACC = \frac{TP + TN}{n}.$$

Ostalih meril (senzitivnost, specifičnost, preciznost), ki izhajajo iz matrike zmot, nismo uporabili, ker imamo v učni množici enako število pozitivnih in negativnih primerov.

Točnost napovedovanja smo preverili tako, da smo učne podatke z metodo prečnega preverjanja na k delih (angl. *k-fold cross validation*) razdelili na $k = 5$ enako velikih delov. Preverjanje poteka v k iteracijah. V vsaki iteraciji množico primerov razdelimo na učno množico, ki predstavlja $k - 1$

delov celotne množice, ter testno množico, ki predstavlja preostali del celotne množice. Model naučimo na $k - 1$ delih učne množice, nato pa ocenimo klasifikacijsko točnost napovedi na preostalih primerih - testni množici. Postopek ponovimo tako, da je vsak od k delov enkrat izločen iz učne množice in izbran za testno množico. Končna klasifikacijska točnost je povprečje klasifikacijskih točnosti posameznih iteracij.

5.2 Rezultati

V tabeli 5.2 so zbrani rezultati prečnega preverjanja. Pri vseh preizkusih smo za napovedovanje uporabili metodo logistične regresije. Kot že omenjeno, je bil namen tega testiranja zgolj preveriti kako dobro se obnesejo implementirane metode predstavitve besedil na enostavnejšem problemu napovedovanja sentimenta, ki je lažji od napovedovanje gibanja tečajev na podlagi besedil. Za čisto osnovo smo vzeli predstavitev besedil z vrečo besed in vrečo nizov. Različna izbira dolžine znakov (pri vreči nizov) oz. števila besed (pri vreči besed) ni bistveno vplivala na klasifikacijsko točnost.

Preizkusili smo tudi napovedovanje z uporabo word2vec vektorjev besed, ki so opisani v poglavju 4.1.3. Implementirali in preizkusili smo obe predlagani metodi za združevanje vektorjev besed v en sam vektor tvita (z uporabo povprečja in z uporabo skupin). Uporabili smo word2vec vektorje besed, naučene na pridobljeni množici tvitov - vzorcu tvitov iz dveh mesecev, ki vsebuje približno 260 milijonov tvitov. Množica ne vsebuje tvitov, ki se uporabljajo za učne podatke ob napovedovanju sentimenta. Vsaka beseda v slovarju metode word2vec je predstavljena s 300-dimenzionalnim vektorjem, celoten slovar pa vsebuje 50000 besed. Klasifikacijska točnost pri uporabi word2vec vektorjev besed se zniža v primerjavi z osnovnimi metodami, ki uporabljajo vrečo besed.

Možen vzrok za tako znižanje točnosti je nezmožnost transformacije besede v word2vec vektor. To se zgodi v primeru, ko besede ni v slovarju metode word2vec, kar pomeni da beseda ni bila prisotna v fazi učenja word2vec vek-

torjev. Tudi avtorji metode word2vec v [20] opozarjajo, da je potrebno za uspešno učenje uporabiti veliko množico besedil. Ta opažanja smo potrdili tako, da smo prešteli delež besed, ki ne morejo biti transformirane v vektorje, ker niso v slovarju metode word2vec. Takih je 27% besed, ki se pojavijo v tvitih. Po naših ocenah bi bilo za doseganje boljše točnosti z uporabo word2vec vektorjev besed zmanjšati delež teh besed. Ocenjujemo, da bi s pomočjo večje množice tvitov lahko dosegli manj kot 10% besed, ki niso prisotne v word2vec slovarju, vendar podrobnejšo analizo v tej smeri prepuščamo za nadaljnje delo.

Preizkušali smo tudi uporabo že izdelanih word2vec vektorjev besed, ki so naučeni na novicah projekta Google News¹. Skupno je učna množica v tem primeru vsebovala 300 milijard besed, v slovarju metode word2vec pa je tri milijone različnih besed. Z uporabo teh vektorjev besed smo dosegali še večji delež besed, ki niso prisotne v slovarju, zato je v nadaljevanju nismo več uporabili.

Zadnje v tabeli rezultatov 5.2 so metode sprotnega izbora atributov, ki so opisane v poglavju 4.2.2. Te metode za predstavitev besedil uporabljajo vrečo besed. Od običajnih metod, ki uporabljajo vrečo besed, se razlikujejo v napovednem modelu. Običajne metode naredijo en napovedni model in uporabijo vnaprej določeno število atributov, metoda sprotnega izbora atributov pa za vsako napoved zgradi nov napovedni model samo z uporabo atributov, ki so prisotni v primeru, za katerega delamo napoved. Iz tabele je razvidno, da metoda sprotnega izbora atributov na uporabljeni množici označenih tvitov izboljša klasifikacijsko točnost v primerjavi z osnovno metodo, ki attribute predstavi z vrečo besed.

Po podatkih iz članka [28] se ljudje glede sentimenta besedila strinjamo le v do 82% primerih. Z rezultatov lahko torej sklepamo, da smo blizu teoretični meji napovedljivosti sentimenta.

Avtorji množice podatkov v članku [19] poročajo o podobnih, sicer rahlo boljših rezultatih klasifikacijske točnosti. Njihov najboljši klasifikator

¹<https://code.google.com/p/word2vec/>

doseže klasifikacijsko točnost 0.83. Vzrok za boljšo točnost njihovih modelov je najverjetneje uporaba dodatnega slovarja sentimentalno označenih besed, ki ga v našem primeru nismo želeli uporabiti. Namen naših eksperimentov na problemu napovedovanja sentimenta je bil preizkusiti metode, ki jih bomo uporabili tudi pri napovedovanju gibanja tečajev, kjer pa ne bi mogli uporabiti slovarja sentimentalno označenih besed.

5.2.1 Regularizacija logistične regresije

Preverili smo, kako regularizacija metode logistične regresije vpliva na klasifikacijsko točnost napovedovanja. Za vsak pristop smo preizkusili tri različne tipe regularizacije, $l1$, $l2$ in *elasticnet*, ter večje število parametrov *alpha* med 1 in $1e-7$. Tip regularizacije je v knjižnici *scikit-learn* določen s parametrom *penalty*.

V tabeli 5.2 posamezna vrstica prikazuje najboljši parameter *penalty* in parameter *alpha* skupaj s klasifikacijsko točnostjo *ACC* za vsak implementiran pristop predstavitve besedil. Pri metodah vreče besed in vreče nizov se bolje obnese uporaba $l1$ regularizacije, pri metodah, ki uporabljajo word2vec vektorje, pa je boljša $l2$ regularizacija. V nadaljevanju smo pri uporabi posameznega pristopa uporabili parametre regularizacije logistične regresije, določene s tem testiranjem.

pristop	opombe	ACC	$penalty$	$alpha$
Vreča besed	-	0.788	l2	1
Vreča besed	TFIDF	0.789	l1	1e-6
Vreča nizov (2, 6)	TFIDF	0.799	l1	1e-6
Vreča nizov (4, 10)	TFIDF	0.794	l1	1e-6
Word2vec (povprečje)	-	0.754	l2	1
Word2vec (povprečje)	TFIDF	0.746	l1	1e-5
Word2vec (skupine)	k=300	0.707	l2	1e-3
Word2vec (skupine)	k=500	0.718	l2	1e-3
Word2vec (skupine)	k=1500	0.733	l2	1e-3
Sprotno, Vreča besed	-	0.787	l2	1
Sprotno, Vreča besed	TFIDF	0.801	l2	1

Tabela 5.2: Tabela prikazuje klasifikacijsko točnost napovedovanja senti-
menta tvitov, ACC , pri uporabi različnih pristopov predstavitve besedil.
Tip regularizacije $penalty$ in faktor $alpha$ sta optimalno določena parametra
logistične regresije s pomočjo prečnega preverjanja.

Poglavje 6

Simulacija trgovanja

Ocenjevanje uspešnosti implementiranih metod pri trgovanju smo izvedli s pomočjo lastnega simulacijskega okolja. Testirali smo različne kombinacije predstavitev besedil in tehničnih podatkov. Primerjali smo tudi, kako izbor razreda, ki ga napovedujemo, vpliva na uspešnost trgovanja oz. napovedovanja gibanja tečaja EURUSD.

Odpiranje posla pomeni virtualni nakup ali prodajo določene količine denarnih sredstev. V resnici se celoten postopek izvede elektronsko, v teoriji pa nakup pri tečaju EURUSD pomeni nakup EUR z USD, prodaja pa pomeni nakup USD z EUR pri upoštevanem menjalnem tečaju v danem trenutku. Pazljivi moramo biti na razliko med nakupno in prodajno ceno, ki je trenutno 0.0003 vrednosti tečaja EURUSD, kar mora upoštevati tudi simulator ob odprtju posla. Če na primer odpremo nakupni posel pri vrednosti tečaja 1.1000, z razliko med nakupno in prodajno ceno 0.0003, lahko posel takoj po odprtju zapremo pri vrednosti tečaja 1.0997, z izgubo 3 točk.

Razlika med simulacijo trgovanja in teoretičnim ocenjevanjem napovedi, ki ga uporablja veliko raziskovalcev s tega področja, je lahko velika. S prečnim preverjanjem lahko dobimo dobro klasifikacijsko točnost, uspešnost takega napovednega modela v resničnem trgovanju pa je lahko zelo slaba. Do tega pride v primeru, ko vrednost tečaja med $open_i$ in $close_i$ v začetku dneva narašča, nato pa v kratkem časovnem obdobju (lahko tudi nekaj minut)

drastično izgubi vrednost, ter preostanek dneva počasi raste, do večera nadoknadi izgubo in zaključi dan s pozitivno razliko v vrednosti $close_i - open_i$. Pri resničnem trgovanju imamo v takem primeru najverjetneje izgubo, saj lahko dosežemo mejne vrednosti izgube, kjer nam banka samodejno zapre posel z izgubo, kljub teoretično pravilni napovedi smeri ob koncu dneva in pozitivni dnevni spremembi $close_i - open_i$. Za razliko od ostalih smo želeli oceniti uspešnost trgovanja v resničnih razmerah, zato smo izdelali simulacijsko okolje za trgovanje, ki posnema vse značilnosti trgovanja, ki jih s prečnim preverjanjem ne moremo upoštevati.

V nadaljevanju je opisan sistem za simulacijo trgovanja in njegovo delovanje. Nato so opisani eksperimenti, ki smo jih izvedli, da bi ugotovili kakšne metode in parametre izbrati za uspešno trgovanje. Preizkusili smo oba tipa metod, metode temeljne in metode tehnične analize. Proti koncu poglavja predstavimo združevanje napovedi različnih metod in opišemo nekaj najboljših rezultatov eksperimentov. Pri vseh simulacijah, razen kjer je opisano drugače, smo za metodo strojnega učenja uporabili logistično regresijo. Uporabili smo tudi regularizacijo logistične regresije, s privzetimi parametri $penalty = l2$ in $alpha = 1$. Dodatno smo preizkusili tudi parametre, ki smo jih določili za posamezen tip predstavitve besedila s prečnim preverjanjem v prejšnjem poglavju.

6.1 Sistem za simulacijo trgovanja

Za potrebe ocenjevanja uspešnosti implementiranih metod smo izdelali sistem, ki simulira pravo trgovanje na valutnem trgu (bolj podrobno opisano v poglavju 2). S simulacijo smo se omejili na obdobje med novembrom 2010 in junijem 2015, kar predstavlja štiri leta in osem mesecev. Razlog za omejevanje obdobja je v tem, da je težko pridobiti točne podatke o vrednosti valutnih tečajev, praktično nemogoče pa je pridobiti tvite izbranih uporabnikov iz obdobja več kot pet let v preteklost.

Zasnova sistema za simulacijo oz. simulatorja je prikazana na sliki 6.1.

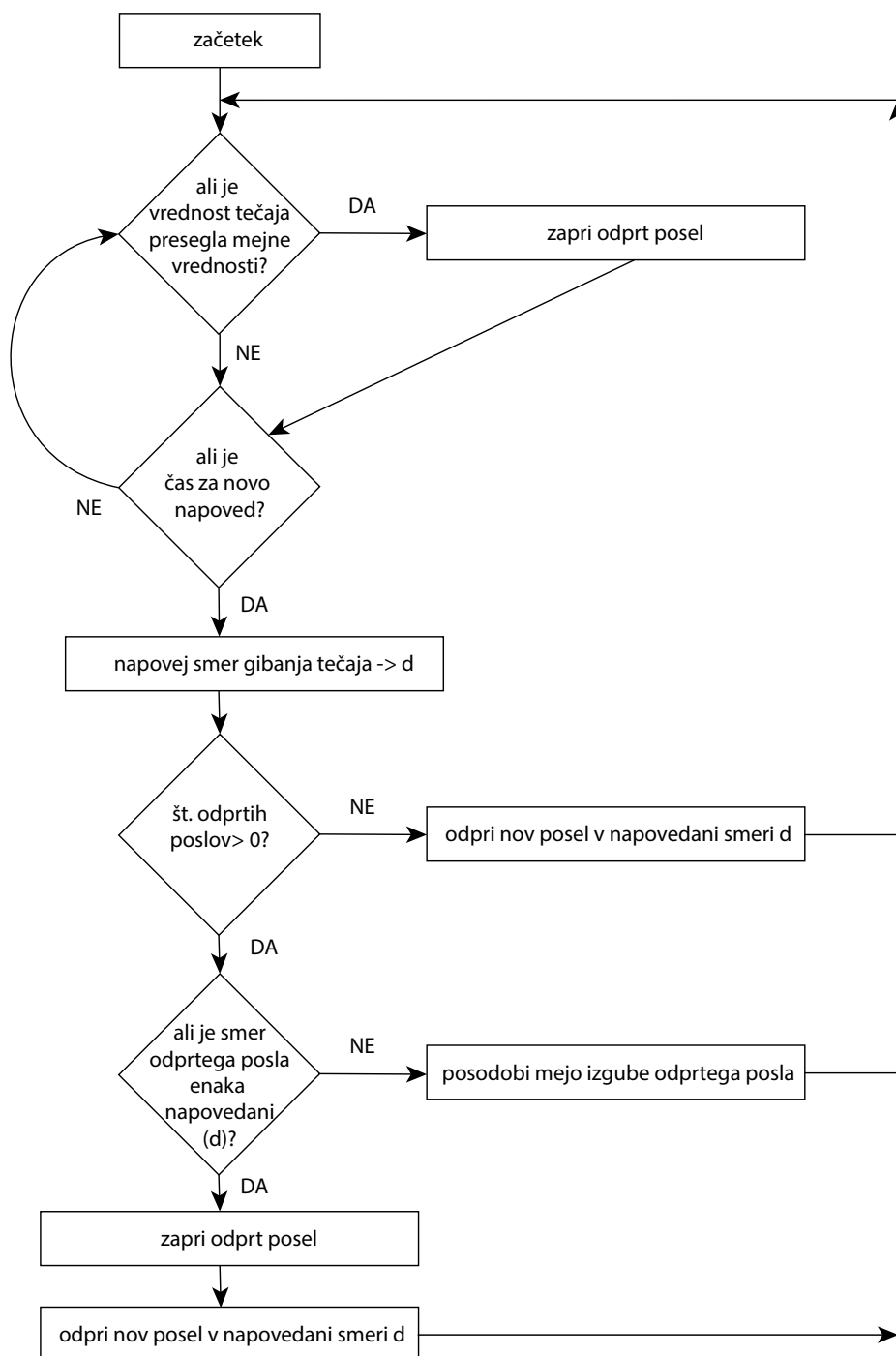
Zasnovan je tako, da omogoča odpiranje in zapiranje poslov glede na napoved napovednega modela. Napoved za odpiranje poslov sistem lahko naredi enkrat dnevno in sicer ob začetku vsakega dneva. Ob primernem času simulator uporabi napovedni model za izdelavo napovedi. Napovedni model torej v začetku dneva t_i glede na izbrane učne podatke iz preteklih dni t_j , $j < i$, izdelava napoved $\hat{d}(t_i) \in \{1, 0, -1\}$.

V primeru napovedi $\hat{d}(t_i) = \pm 1$ simulator preveri odprte posle. Omejitev razvitega simulatorja je en odprt posel naenkrat. Če ni odprtih poslov, odpre nov posel v napovedani smeri (± 1). Če že obstaja odprt posel, preveri njegovo smer $\hat{d}(t_{i-x})$, kjer je t_{i-x} dan v preteklosti, ko je bil posel odprt. V primeru iste smeri $\hat{d}(t_{i-x}) = d(t_i)$ sistem zgolj posodobi mejno vrednost za zaustavitev izgube, v primeru nasprotnih vrednosti $\hat{d}(t_{i-x}) = -\hat{d}(t_i)$ pa posel, odprt na začetku dneva t_{i-x} , zapre in nato odpre nov posel v napovedani smeri gibanja valutnega tečaja $\hat{y}(t_i)$.

V tej nalogi smo se omejili na odpiranje posla enkrat na dan in sicer ob začetku dneva (po časovnem pasu GMT). Naloga simulatorja je preverjanje, če je kateri od odprtih poslov presegel mejne vrednosti, ki so določene pri odpiranju posla. To preverjanje se mora izvajati v manjšem časovnem intervalu, saj je pogosto, da vrednost tečaja znotraj enega dneva močno niha, ob koncu dneva pa se ustali na vrednosti blizu začetni vrednosti istega dne. V ta namen simulator vsako polno uro preveri mejne vrednosti in v primeru presežene vrednosti zapre odprt posel.

6.1.1 Primer delovanja simulatorja

Ob začetku simulacije določimo začetno stanje sredstev. Za vse eksperimente smo določili začetno stanje sredstev 1000 USD. Ob vsakem zapiranju posla se stanje sredstev posodobi za nastali dobiček ali izgubo pravkar zaprtega posla. V tabeli 6.1 je prikazan primer izpisa trgovanja simulatorja. Vrstica 1 tabele 6.1 predstavlja odpiranje posla 22. maja ob polnoči v negativni smeri. Posel je odprt pri vrednosti 1.2809, mejna vrednost izgube je nastavljena na vrednost 1.2834. Naslednji dan ob isti uri simulator na podlagi ravnokar za-



Slika 6.1: Diagram, ki prikazuje delovanje simulacijskega okolja za trgovanje

	Čas	Napoved	Akcija	Odperto	Meja	Zaprto	Stanje
1	22.5. 00:00	-1	odpri (prodaj)	1.2809	1.2834		1000.00
2	23.5. 00:00	-1	posodobi	1.2809	1.2870		1000.00
3	24.5. 00:00	-1	posodobi	1.2809	1.2906		1000.00
4	25.5. 00:00	-1	posodobi	1.2809	1.2893		1000.00
5	28.5. 00:00	-1	zapri (obratno)	1.2809	1.2893	1.2566	1038.69
6	28.5. 00:00	1	odpri (kupi)	1.2569	1.2508		1038.69
7	29.5. 00:00	1	posodobi	1.2569	1.2527		1038.69
8	29.5. 01:00	1	zapri (izguba)	1.25699	1.2527	1.2527	1032.08

Tabela 6.1: Primer izpisa delovanja simulatorja

ključenega dneva in preteklih dni ponovno izdela napoved v negativni smeri. Ker je en posel v isti smeri že odprt, simulator zgolj posodobi mejno vrednost izgube na 1.2870. Podobno se zgodi še v naslednjih dveh dneh. Simulator 28. maja izdela napoved v pozitivni smeri. Ker je en posel v nasprotni smeri trenutne napovedi že odprt, ga zapre pri vrednosti tečaja 1.2566, ter odpre novega v napovedani smeri pri vrednosti tečaja 1.2569 z mejno vrednostjo izgube pri vrednosti tečaja 1.2508. Naslednji dan simulator ponovno napove gibanje tečaja v pozitivni smeri, zato zgolj posodobi mejno vrednost izgube na vrednost tečaja 1.2527. Eno uro za tem vrednost tečaja pade pod mejno vrednost izgube 1.2527, zato simulator posel zapre z izgubo.

Zapiranje posla v vrstici 5 tabele 6.1 je pridelalo 243 točk razlike. Dobiček je odvisen od velikosti posla. V simulaciji uporabljamo 1000 enot, dobiček v denarnih sredstvih pa izračunamo glede na zaključno vrednost posla in razliko v točkah, $0.0243 \cdot 1.2566 \cdot 1000 \cdot 1.2566 = 38.37$. Upoštevati moramo še obresti na posel s pozitivnim stanjem, zato prištejemo še 32 centov pozitivnih obresti in dobimo končni dobiček posla 38.69 USD.

6.1.2 Naključno trgovanje

Pred resnim testiranjem implementiranih metod in njihovih parametrov smo najprej testirali uspešnost trgovanja z naključnim klasifikatorjem, ter tako preverili delovanje simulatorja.

Simulacijo trgovanja smo ponovili 100-krat. Povprečna vrednost končnega stanja pri 100 ponovitvah je 512 USD, kar je skoraj polovica manj kot vrednost začetnih sredstev in predstavlja izgubo. Dobiček oziroma končno vrednost sredstev nad 1000 USD je doseglo le 9 od 100 simulacij naključnega trgovanja, kar je bilo pričakovano.

6.2 Izbira primerne razreda

Izbira razreda je lahko ključnega pomena pri uspešnosti napovedovanja. V začetku tega poglavja smo opisali način preverjanja uspešnosti, ki v teoriji (s prečnim preverjanjem) lahko doseže veliko točnost, v pravih okoliščinah med trgovanjem pa deluje slabo. Lahko bi tudi trdili, da je uspešnost odvisna od izbire razreda.

Določili smo dva razreda, c_i^{CHANGE} in c_i^{ATR} . Razred c_i^{CHANGE} predstavlja (pozitivno ali negativno) spremembo vrednosti tečaja v enem dnevu. Kot izboljšavo smo definirali še razred c_i^{ATR} . Njegova vrednost je lahko pozitivna ali negativna, kot v primeru razreda c_i^{CHANGE} . Razred c_i^{ATR} predstavlja tisto smer (pozitivno ali negativno), ki je prva dosegla spremembo v vrednosti tehničnega indikatorja ATR_i^{20} (poglavje 2.2) v začetku dneva t_i . Z drugimi besedami, razred je pozitiven, če je vrednost tečaja od začetka dneva i prej dosegla vrednost $open_i + ATR_{i-1}^{20}$, kot pa $open_i - ATR_{i-1}^{20}$, in obratno za negativen razred.

c_i^{CHANGE} **sprememba cene v zadnjem dnevu** je diskreten razred, definiran z absolutno spremembo vrednosti tečaja v enem dnevu, ki je določena z

$$c_i^{CHANGE_{ABS}} = open_{t_i} - open_{t_{i-1}} .$$

Ker se ukvarjamo s klasifikacijo, smo na preprost način omejili vrednosti na pozitiven in negativen razred s predpisom

$$c_i^{CHANGE} = sign(c_i^{CHANGE_{ABS}}) .$$

c_i^{ATR} kateri **ATR20 bo dosežen prej** je diskreten razred, ki nam pove, katero vrednost ATR_i^{20} , pozitivno ali negativno, je vrednost tečaja dosegla najprej, v obdobju od začetka dneva t_i . Finančni indikator ATR_i^{20} , opisan v poglavju 2.2, predstavlja povprečno 20-dnevno spremembo tečaja in se uporablja kot pričakovana vrednost spremembe v trenutnem dnevu. Razred c_i^{ATR} je določen z

$$c_i^{ATR} = \begin{cases} 1 & \text{če je najprej dosežen } open_{t_i} + ATR_i^{20} \\ 0 & \text{če je najprej dosežen } open_{t_t} - ATR_i^{20} \end{cases}.$$

6.2.1 Eksperimentalni rezultati

S testiranjem na simulatorju trgovanja smo preizkusili kako izbira razreda vpliva na stanje sredstev po končani simulaciji. Seznam najboljših testiranj je priložen v dodatku A. V tabeli 6.2 so predstavljeni rezultati simulacije trgovanja štirih izbranih napovednih modelov, dva sta iz skupine temeljne analize in za vhod uporabljata besedilne attribute, ter dva modela iz skupine tehnične analize (glej 2.2.1), ki uporabljata tehnične attribute.

Modela *tehnični 1* in *tehnični 2* se razlikujeta zgolj v izbranih atributih množice podatkov. Atributi napovednega modela *tehnični 1* so zadnjih pet dnevnih sprememb tečaja, napovedni model *tehnični 2* pa vsebuje tudi attribute tehničnih indikatorjev.

Oba napovedna modela, ki za vhod uporabljata besedilne attribute, z uporabo razreda c_i^{CHANGE} naredita izgubo, medtem ko z uporabo razreda c_i^{ATR} naredita pozitivno razliko od začetnih sredstev v višini 1000 USD. Tehnična modela, ki uporabljata zgolj tehnične attribute, v vsakem primeru naredita izgubo, ki se bistveno ne razlikuje od izbire razreda. Glede na te rezultate je za uspešno trgovanje bolje uporabiti razred c_i^{ATR} . Vse simulacije v nadaljevanju tega poglavja uporabljajo razred c_i^{ATR} .

metoda	c_i^{CHANGE}		c_i^{ATR}	
	sredstva	ACC	sredstva	ACC
vreča nizov	562.19	0.47	2232.32	0.64
vreča besed	431.61	0.47	2107.82	0.64
tehnični 1	734.74	0.36	872.61	0.50
tehnični 2	799.55	0.46	787.71	0.55

Tabela 6.2: Primerjava sredstev po končani simulaciji izbranih napovednih modelov v kombinaciji z definiranimi razredoma

6.3 Primerjava različnih predstavitev besedila

Implementirali smo štiri različne predstavitve besedil, s katerimi smo tvite predstavili v atributni obliki, primerni za metode strojnega učenja. To so vreča besed, vreča nizov, povprečje word2vec in uvrščanje word2vec. Metode so bolj podrobno opisane v poglavju 4.1.

Rezultati simulacije so podani v tabeli 6.3. Opazimo lahko, da so vse metode v obdobju testiranja (približno pet let) podvojile začetno stanje sredstev. Opazimo tudi malenkost višje končno stanje pri uporabi vreče besed in vreče nizov, kot pri word2vec metodah, medtem ko za klasifikacijsko točnost velja ravno obratno in je višja pri word2vec metodah.

Razlika med klasifikacijsko točnostjo in uspešnostjo trgovanja nastane zaradi značilnosti trgovanja. Kot smo že omenili, lahko napovedni model pravilno napove vrednost razreda, vendar to pri resničnem trgovanju ne pomeni zagotovljenega dobička. Do takih razlik pride, ko je trg visoko likviden in prihaja do večjih nihanj vrednosti tečaja.

6.4 Izbira ustreznega števila učnih primerov

Simulator trgovanja ob začetku vsakega dneva s pomočjo izbranega napovednega modela izdelava napoved smeri gibanja tečaja. Za izdelavo napovedi ob začetku dneva t_i vedno na novo naučimo napovedni model na n prete-

metoda	stanje	ACC
vreča nizov	2232.32	0.60
vreča besed	2107.82	0.60
povprečje word2vec	1996.67	0.64
uvrščanje word2vec	1941.07	0.64

Tabela 6.3: Primerjava implementiranih predstavitev besedil glede na sredstva po končani simulaciji in klasifikacijske točnosti

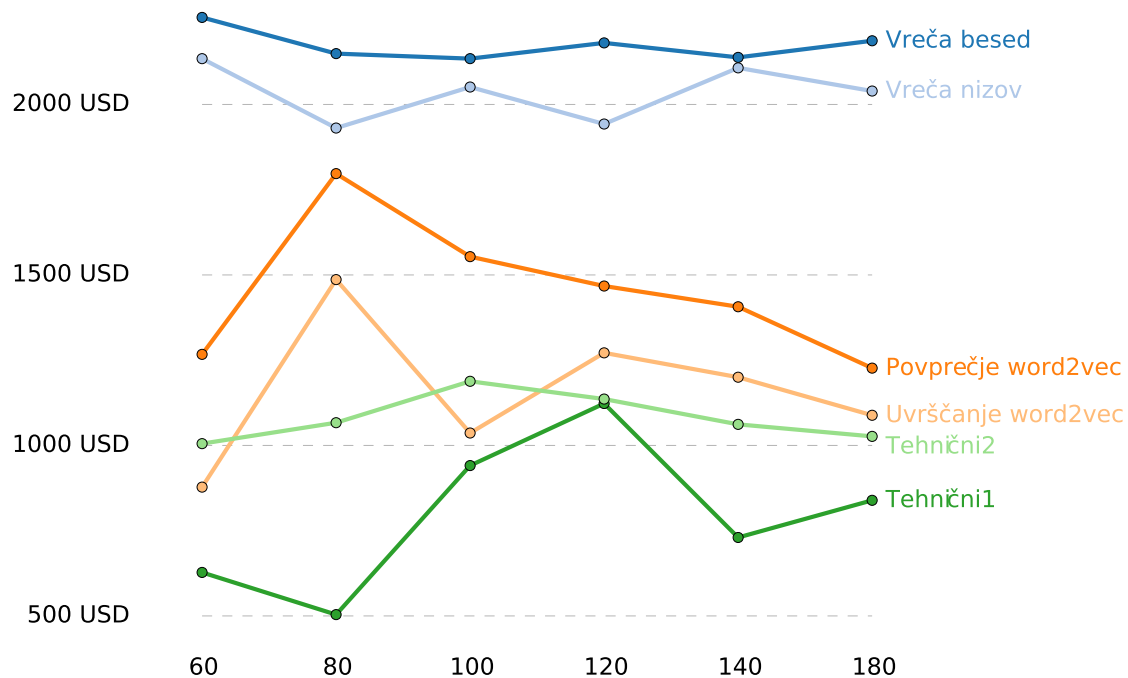
klih primerih. S tem zagotovimo hitro prilagajanje modela na nove trende in hkrati omejimo učno množico na zadnjih n dni. Primer kratkotrajnega trenda je pojav tвитov, ki vsebujejo nekaj novih ključnih besed in so lahko zelo informativni in aktualni v krajšem časovnem intervalu. Če bi v nasprotnem primeru zgradili le en napovedni model, bi teh nekaj ključnih besed, ki so se pojavile v zadnjih nekaj primerih (dneh) zagotovo predstavljalo manjši vpliv v napovednem modelu. Tak napovedni model bi iskal pravila na celotni množici in ne bi zaznal kratkotrajnih manjših trendov.

Podobne pristope smo zasledili med opisi sorodnih del članka [1]. Nekateri raziskovalci so za učenje uporabljali zgolj fiksno dolžino zadnjih 1000 ur ali le primere v zadnje pol leta.

Zanimalo nas je kako izbira števila učnih primerov n vpliva na uspešnost trgovanja, zato smo na simulatorju preizkusili šest napovednih modelov (4 uporabljajo besedilne attribute, 2 modela pa tehnične attribute). Za vsakega od napovednih modelov smo simulirali trgovanje s 60, 80, 100, 120, 140 in 180 učnimi primeri. Slika 6.3 prikazuje rezultate simulacije.

Iz slike 6.3 takoj opazimo, da sta napovedna modela vreča besed in vreča nizov glede na končno vrednost tesno skupaj in boljši od preostalih. Glede na izbor števila učnih parametrov ni opaznejših razlik v uspešnosti teh dveh napovednih modelov.

Nasprotno pri napovednih modelih word2vec opazimo najboljšo uspešnost pri izbiri 80 učnih primerov, z večanjem števila primerov pa se uspešnost samo še slabša. Pri tehničnih napovednih modelih je uspešnost najboljša



Slika 6.2: Primerjava uspešnosti trgovanja različnih metod v kombinaciji z različnim številom učnih primerov

med 100 in 120 učnimi primeri. *tehnični1* je ne glede na izbiro števila učnih primerov naredil izgubo, *tehnični2* pa je za las nad začetno vrednostjo trgovinskih sredstev.

6.5 Izbira ustreznih uteži

Pri izdelavi atributov za predstavitev primera s tehničnimi podatki, še posebno pa z besedili, se navadno omejimo na neko časovno okno. Če izberemo okno enega dneva, so atributi enega primera sestavljeni na podlagi tвитov znotraj izbranega dnevnega intervala.

Časovno okno pri izdelavi atributov je povsem neodvisno od izbranega intervala izdelave napovedi, ki je v tem delu vsak začetek dneva. Attribute lahko izdelamo tudi s poljubno večjim časovnim oknom, še vedno pa napovedujemo gibanje vsak začetek dneva. V primeru da izberemo časovno okno za izdelavo atributov tri dni, za napoved gibanja ob začetku dneva t_i uporabimo tvite med vključno t_{i-3} in t_{i-1} , kar pomeni, da uporabimo tvite preteklih treh dni.

Za izdelavo atributov smo implementirali funkcijo, ki sprejme seznam uteži w . Posamezna utež w_j h končnim atributom primera za dan t_i doda attribute dneva t_{i-j} utežene z w_j . Nabor uteži, ki smo jih ročno izbrali in preizkusili na šestih napovednih modelih, je podan v tabeli 6.4.

Če na primer uporabimo nabor uteži $u^{(4)} = [1.0, 0.8, 0.6]$, so atributi primera x_i , ki predstavlja tvite do vključno dneva t_i , izračunani s predpisom

$$x_i = u_0^{(4)} x_i + u_1^{(4)} x_{i-1} + u_2^{(4)} x_{i-2} = 1.0x_i + 0.8x_{i-1} + 0.6x_{i-2} .$$

Na sliki 6.3 je prikazana uspešnost z uporabo različnih uteži za izdelavo atributov. Iz slike je jasno razvidno, da v večini uspešnost trgovanja med $u^{(1)}$ do $u^{(3)}$ narašča, kar pomeni da so bolj uspešni napovedni modeli, ki pri izdelavi atributov upoštevajo tudi več preteklih dni.

Med utežmi $u^{(3)}$ do $u^{(6)}$, ki upoštevajo pretekle tri dni, je pri vreči besed in vreči nizov najboljša izbira $u^{(3)}$. Pričakovali smo ravno obraten rezultat,

oznaka	nabor uteži
$u^{(1)}$	[1.0]
$u^{(2)}$	[1.0, 1.0]
$u^{(3)}$	[1.0, 1.0, 1.0]
$u^{(4)}$	[1.0, 0.8, 0.6]
$u^{(5)}$	[1.0, 0.5, 0.3]
$u^{(6)}$	[0.7, 0.2, 0.1]
$u^{(7)}$	[0.45, 0.2, 0.15, 0.15, 0.05]
$u^{(8)}$	[0.35, 0.2, 0.15, 0.1, 0.1, 0.05, 0.05]
$u^{(9)}$	[1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0]
$u^{(10)}$	[1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0]

Tabela 6.4: Seznam ročno izbranih uteži za uteženo izdelavo atributov

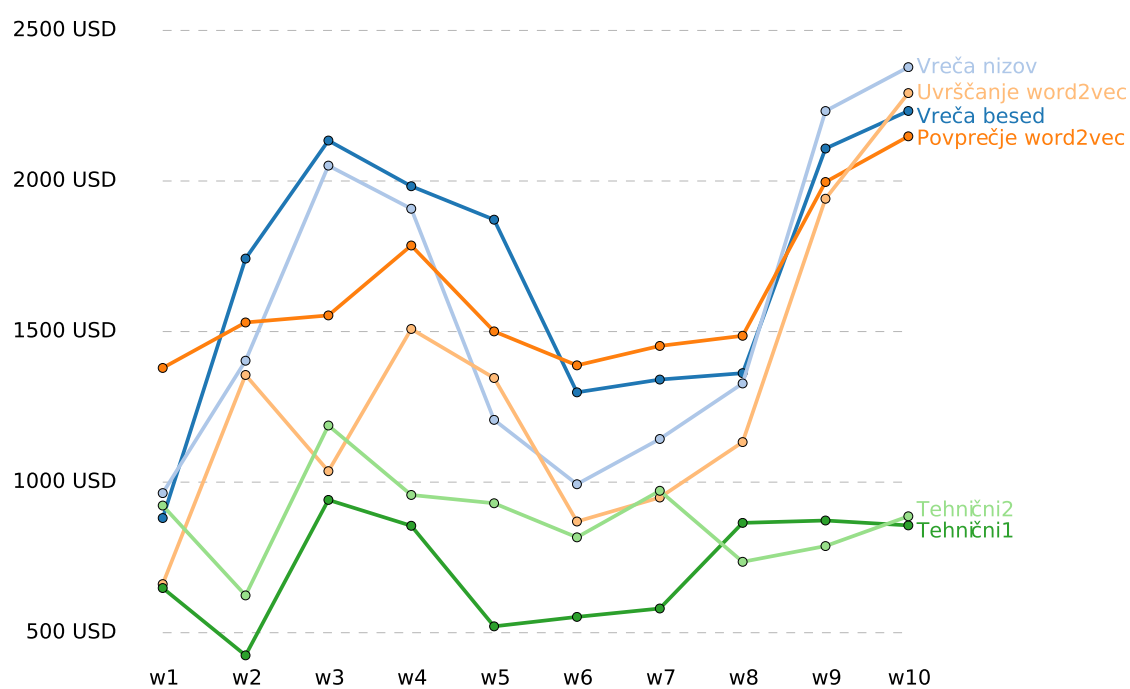
da vpliv atributov na napovedovanje pada s starostjo, vendar to velja zgolj za metodo povprečje word2vec.

Zadnja dva vektorja uteži, $u^{(9)}$ in $u^{(10)}$ imata najvišjo uspešnost trgovanja med štirimi napovednimi modeli, ki so dobičkonosni. Možna razlaga za to je pomanjkanje neničelnih atributov pri uporabi manjšega časovnega okna. Enodnevno časovno okno lahko vsebuje le nekaj deset neničelnih atributov, pri uporabi uteži $u^{(10)}$ pa se lahko število neničelnih atributov poveča do 12-krat, kar lahko pripomore k boljši natančnosti napovedi.

6.6 Združevanje napovedi

Z združevanjem napovedi smo preizkusili uspešnost trgovanja z uporabo tehnične analize in temeljne analize skupaj. Uporabljeni metodi združevanja napovedi, glasovanje in metoda stacking, sta bolj podrobno opisani v poglavju 4.2.3. Z metodami združevanja smo združili napovedi posameznih modelov.

Omejili smo se na pet ročno določenih množic napovednih modelov, ki so podani v tabeli 6.5. Uporabljena metoda strojnega učenja vseh posameznih



Slika 6.3: Primerjava uspešnosti trgovanja glede na izbiro uteži pri izdelavi atributov

modelov je metoda logistične regresije. Množica modelov *top6* vsebuje šest najboljših posamičnih napovednih modelov. Vseh šest uporablja zgolj besedilne attribute. Ročno smo definirali še štiri množice napovednih modelov (*modeli1* - *modeli4*), ki vsebujejo napovedne modele z besedilnimi atributi in modele s tehničnimi atributi. Množica modelov *modeli1* vsebuje samo modele z utežmi zadnjih treh dni, $u^{(3)}$, množice *modeli2*, *modeli3*, *modeli4* pa vsebujejo mešane modele. Množica *modeli3* vsebuje pretežno modele z utežmi dlje v zgodovino, *modeli4* pa tri modele s krajšimi utežmi in tri modele z daljšimi utežmi.

Metoda stacking za svoje delovanje potrebuje učno množico in testno množico primerov. Množico primerov (trgovalnih dni) smo razdelili na dve enako veliki polovici. Na prvi polovici se je metoda stacking s pomočjo metod strojnega učenja naučila pravil združevanja, nato pa smo na drugi polovici primerov izvedli simulacijo trgovanja. Pri primerjavi rezultatov z drugimi napovednimi modeli moramo upoštevati polovični časovni interval trgovanja. Preizkusili smo tri metode strojnega učenja, ki jih uporablja metoda stacking za meta-učenje, naivni Bayes, logistična regresija in metoda naključnih dreves. Rezultati združevanja napovedi so zbrani v tabeli 6.6.

Najboljša izmed množice modelov glede na eksperimentalne rezultate je množica najboljših šestih modelov *top6*. Pri tej množici je najboljša metoda združevanja napovedi metoda stacking z uporabo metode naivni Bayes, zelo blizu pa je tudi kombinacija z metodo logistične regresije. Presenetljivo je pri množici modelov *top6* precej uspešna tudi metoda glasovanja.

Pri *modeli1* je najboljša metoda stacking z logistično regresijo. Izredno slabo se odreže metoda stacking z uporabo naključnih gozdov. Med množicami modelom *modeli1* - *modeli4* je najboljša *modeli2*, ki vsebuje različne napovedne modele z različnimi predstavitvami besedil.

Zanimivo pri množici modelov *modeli4* je splošno slabša uspešnost v primerjavi z ostalimi, izstopa pa relativno uspešna kombinacija z metodo logistične regresije.

Najboljša metoda združevanja napovedi je metoda stacking v kombinaciji

množica	predstavitev atr.	učnih primerov	uteži
top6	vreča nizov	100	$u^{(10)}$
	vreča nizov	120	$u^{(10)}$
	vreča nizov	120	$u^{(9)}$
	vreča nizov	80	$u^{(9)}$
	povprečje w2v	140	$u^{(10)}$
	povprečje w2v	60	$u^{(10)}$
modeli1	vreča besed	80	$u^{(3)}$
	vreča nizov	80	$u^{(3)}$
	vreča nizov	100	$u^{(3)}$
	vreča nizov	120	$u^{(3)}$
	povprečje word2vec	140	$u^{(3)}$
	Skupine word2vec	80	$u^{(3)}$
	vreča besed (izbor atr.)	100	$u^{(3)}$
	tehnični1	100	$u^{(3)}$
	tehnični2	100	$u^{(3)}$
	tehnični2	180	$u^{(3)}$
modeli2	vreča nizov	80	$u^{(3)}$
	vreča besed	120	$u^{(3)}$
	povprečje word2vec	140	$u^{(9)}$
	tehnični1	100	$u^{(9)}$
	tehnični2	180	$u^{(9)}$
modeli3	vreča nizov	80	$u^{(10)}$
	vreča nizov	180	$u^{(10)}$
	povprečje word2vec	140	$u^{(9)}$
	tehnični1	80	$u^{(9)}$
	tehnični1	100	$u^{(3)}$
	tehnični2	100	$u^{(9)}$
	tehnični2	180	$u^{(3)}$
modeli4	vreča nizov	80	$u^{(10)}$
	vreča besed	120	$u^{(3)}$
	Skupine word2vec	80	$u^{(3)}$
	tehnični1	100	$u^{(10)}$
	tehnični2	100	$u^{(10)}$
	tehnični2	180	$u^{(3)}$

Tabela 6.5: Tabela prikazuje ročno definirane množice modelov za testiranje metod za združevanje napovedi. Uporabljena metoda strojnega učenja vseh posameznih modelov je logistična regresija, poleg vsakega pa so zapisani ostali uporabljeni parametri: tehnika predstavitve besedila, število učnih primerov in uporabljene uteži.

	glasovanje	metoda stacking		
množica		logistična reg.	naivni Bayes	naključni gozd
top6	1405.45	1889.3	1929.92	328.04
modeli1	1316.91	1373.28	1271.44	833.18
modeli2	1224.31	1602.47	1415.07	1028.19
modeli3	1205.03	1337.34	1187.98	1392.96
modeli4	1057.07	1507.75	806.91	401.28

Tabela 6.6: Rezultati simulacije trgovanja z metodami združevanja napovedi

z logistično regresijo. Rahlo slabše (razen pri množici modelov *top6*) je metoda stacking z metodo naivni Bayes, zelo slabo pri združevanju pa se izkaže metoda naključnih gozdov. Metoda glasovanja v splošnem deluje rahlo slabše kot najboljše metode, vendar je glede na enostavnost zelo uspešna.

Pri primerjanju rezultatov z ostalimi metodami moramo upoštevati polovično obdobje trgovanja, to je dve leti in štiri mesece. Če preprosto sklepamo, da bi metode na celotnem obdobju trgovanja dosegale isto uspešnost oz. dvojno vrednost sredstev, kot na polovičnem trgovanju, bi lahko presegle dobiček v višini 2800 USD. Tako lahko sklepamo, da so metode združevanja boljše od posamičnih napovednih modelov.

6.7 Najboljših 10 napovednih modelov

Skupno smo preizkusili več kot 1000 napovednih modelov. V dodatku A je skrajšan seznam preizkušenih kombinacij parametrov s pripadajočo uspešnostjo napovednih modelov. Parametri desetih najbolj dobičkonosnih napovednih modelov so predstavljeni v tabeli 6.7. Na sliki 6.4 je prikazano stanje sredstev med trgovanjem za napovedna modela, ki uporabljata vrečo nizov (uporabljenih 100 učnih primerov) in povprečje word2vec (uporabljenih 120 učnih primerov). Najboljših 10 modelov uporablja logistično regresijo s privzetimi parametri regularizacije ($penalty = l2$ in $alpha = 1$). Modeli s spremenjenimi

parametri so prav tako podani v dodatku A. Glede na uspešnost teh modelov lahko rečemo, da izbira parametrov regularizacije ni odločilno vplivala na uspešnost simulacije trgovanja.

Rezultati simulacij kažejo, da so napovedni modeli temeljne analize boljši od modelov tehnične analize, ki povečini delajo izgubo ali zelo majhen dobiček. Zanimivo ni večjih razlik med uporabljenimi metodami predstavitve besedil, saj so precej skupaj po uspešnosti trgovanja, vendar le z optimalnim izborom atributov, ki smo jih pridobili z izvajanjem simulacij. Dejavnik, ki je najbolj vplival na uspešnost modela, je, glede na rezultate, izbira uteži pri izdelavi atributov. Večina najboljših modelov je namreč takih z utežmi $u^{(9)}$ ali $u^{(10)}$, kar pomeni, da so atributi primera sestavljeni iz preteklih sedem ali dvanaest dni tvtov oz. vrednosti tehničnih indikatorjev (odvisno od tipa metode).

Iz slike 6.4 je razvidno konstantno naraščanje sredstev brez večjih negativnih intervalov. Če stanje sredstev primerjamo z gibanjem tečaja EURUSD (na sliki 3.2) opazimo nekaj pozitivnih in nekaj negativnih obdobj gibanja tečaja. Razviti napovedni modeli napovedujejo dobro, ne glede na trend gibanja tečaja, kar je pri trgovanju ključno.

Presenečeni smo nad raznolikostjo najboljših napovednih modelov, še posebno nad dejstvom, da je med najboljšimi prav vsaka od implementiranih predstavitev besedil: vreča nizov, vreča besed in predstavitev z word2vec vektorji. Pričakovali smo prevlado predstavitev besedil z word2vec vektorji besed, vendar so zgolj primerljive z enostavnejšima, vrečo besed in vrečo nizov.

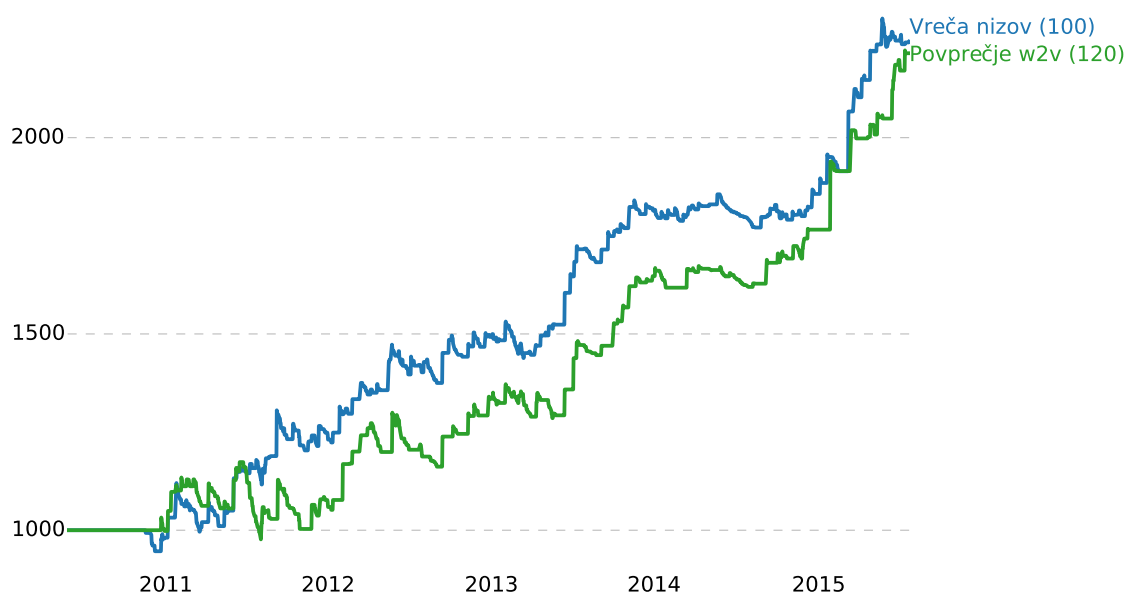
6.8 Primerjava s trgovanjem na valutnem trgu

Po podatkih zbranih iz spletne strani BarclayHedge¹, najboljši investicijski skladi v povprečju zaključijo leto s 40% dobičkom, kar v petletnem obdobju predstavlja več kot 300% dobička.

¹<https://www.dukascopy.com/fxcomm/fx-article-contest/?How-Big-Guys-Perform-In=&action=read&id=590&language=en>

model	učnih primerov	uteži		sredstva
vreča nizov	100	$u^{(10)}$		2378.02
vreča nizov	120	$u^{(9)}$		2374.19
uvrščanje word2vec	60	$u^{(10)}$		2365.65
povprečje word2vec	140	$u^{(10)}$		2359.45
vreča besed	60	$u^{(4)}$		2330.67
vreča besed	140	$u^{(10)}$		2329.69
povprečje word2vec	120	$u^{(9)}$		2327.27
vreča nizov	80	$u^{(9)}$	da	2312.84
vreča nizov	140	$u^{(9)}$	da	2304.92
uvrščanje word2vec	100	$u^{(10)}$		2291.95

Tabela 6.7: 10 najboljših napovednih modelov



Slika 6.4: Prikaz gibanja vrednosti sredstev med simulacijo dveh izmed najboljših napovednih modelov

Po drugi strani študija [29] poroča o zelo slabi uspešnosti trgovanja posameznikov, ki trgujejo na valutnem trgu. Po njihovih podatkih 89% trgovalcev v obdobju štirih let izgubi vsa trgovalna sredstva.

Če primerjamo uspešnost trgovanja najboljših simulacij z uspešnostjo investicijskih skladov, so slednji daleč bolj uspešni. V primerjavi z uspešnostjo trgovanja posameznikov lahko rečemo, da razviti modeli trgujejo precej bolje, saj najboljši v petletnem obdobju naredijo več kot 100% dobička.

Poglavje 7

Sklepne ugotovitve

Namen tega dela je bil preizkusiti že znane metode strojnega učenja in predstavitev besedil na problemu napovedovanja vrednosti valutnih tečajev oz. oceniti njihovo uspešnost pri trgovanju.

Implementirane metode smo preizkusili na manjšem in enostavnejšem problemu - klasifikaciji tvitov v pozitiven ali negativen razred, ki predstavlja sentiment. Ugotovili smo, da delujejo zadovoljivo že brez večjega truda pri izbiri parametrov.

V simulacijskem okolju za trgovanje smo preizkusili in ovrednotili različne predstavitve besedil, metode strojnega učenja, njihove parametre in parametre trgovanja simulacijskega okolja. Z eksperimenti, predstavljenimi v poglavju 6.2, smo ugotovili, da na uspešnost napovedovanja vpliva izbira klasifikacijskega razreda, ki pripada primeru.

Množica tvitov, uporabljena za napovedovanje, vsebuje zgolj 1.2 milijona tvitov, oz. povprečno 28.3 tvita na dan. Z upoštevanjem uteži pri grajenju atributov (poglavje 6.5) en primer predstavimo s tviti večih dni, ter tako zberemo dovolj neničelnih atributov, da so uporabni pri napovedovanju. V prihodnje nameravamo učno množico tvitov razširiti in preveriti kako večja množica izboljša uspešnost trgovanja, ter ponovno preveriti kako upoštevanje tvitov iz preteklosti vpliva na uspešnost trgovanja z uporabo večje množice tvitov.

Implementirali smo tudi dodatno razširitev metode vreče besed - sproten izbor atributov [1], vendar se glede na rezultate, v nasprotju z avtorji, ni izkazala za ključno izboljšavo.

Skupno smo preizkusili več kot 1000 različnih napovednih modelov, najboljših 600 glede na končno stanje sredstev je podanih v dodatku A. Presenetljivo so med desetimi najbolj uspešnimi vse implementirane metode predstavitev besedil. Predstavitve besedil z word2vec vektorji besed imajo po našem mnenju še veliko možnosti za izboljšave. Prva je učenje word2vec vektorjev na veliko večji množici tvitov in s tem tudi zmanjšanje deleža besed, ki niso v slovarju word2vec, ki jih sedaj preprosto ignoriramo, ker jih ne moremo predstaviti z ustreznim atributom. Druga izboljšava je uporaba tretjega načina združevanja vektorjev besed v vektorje tvitov. Avtor knjižnice Gensim je implementiral svojo različico transformacije dokumentov v en sam vektor, vendar je nismo uspeli uporabiti in jo prepuščamo za nadaljnje raziskovanje.

Glede na rezultate simulacij so metode temeljne analize boljše od metod tehnične analize. Združevanje napovedi se je v skladu s pričakovanji izkazalo za dobro izboljšavo uspešnosti posamičnih napovednih modelov (rezultati v poglavju 6.6). Najboljša metoda za združevanje je po naših testih metoda stacking z uporabo metode logistične regresije ali metode naivni Bayes.

Simulacije trgovanja smo preizkusili zgolj na valutnem paru EURUSD. Ta valutni par predstavlja eno četrtno vsega trgovanja na valutnem trgu, zato smo ga izbrali za raziskovanje. V prihodnje želimo simulacije razširiti tudi na druge večje valutne pare, kot sta USDJPY in GBPUSD. V magistrskem delu smo pokazali, da je Twitter primeren vir informacij, ki dovolj dobro opišejo finančno dogajanje za uspešno trgovanje, kar smo pokazali z velikim številom uspešnih simulacij trgovanja.

Uspešnost razvitih metod je slabša od uspešnosti najboljših investicijskih skladov, ki dosegajo več kot 300% dobičke v petletnem obdobju. V primerjavi s trgovanjem posameznikov, ki po podatkih [29] v istem obdobju v 89% primerih izgubijo vsa sredstva, so razvite metode trgovanja z uporabo optimalno izbranih parametrov bolj uspešne, saj glede na simulacije v petletnem

obdobju podvojijo vložena sredstva.

Literatura

- [1] A. K. Nassirtoussi, S. Aghabozorgi, T. Y. Wah, D. C. L. Ngo, Text mining of news-headlines for FOREX market prediction: A multi-layer dimension reduction algorithm with semantics and sentiment, *Expert Systems with Applications* 42 (1) (2015) 306 – 324. doi:<http://dx.doi.org/10.1016/j.eswa.2014.08.004>.
URL <http://www.sciencedirect.com/science/article/pii/S0957417414004801>
- [2] J. Bollen, H. Mao, X. Zeng, Twitter mood predicts the stock market, *Journal of Computational Science* 2 (1) (2011) 1 – 8. doi:<http://dx.doi.org/10.1016/j.jocs.2010.12.007>.
URL <http://www.sciencedirect.com/science/article/pii/S187775031100007X>
- [3] M. Kaya, M. Karşligil, Stock price prediction using financial news articles, in: *Information and Financial Engineering (ICIFE)*, 2010 2nd IEEE International Conference on, 2010, pp. 478–482. doi:[10.1109/ICIFE.2010.5609404](https://doi.org/10.1109/ICIFE.2010.5609404).
- [4] T. Rao, S. Srivastava, Tweetsmart: Hedging in markets through twitter, in: *Emerging Applications of Information Technology (EAIT)*, 2012 Third International Conference on, 2012, pp. 193–196. doi:[10.1109/EAIT.2012.6407894](https://doi.org/10.1109/EAIT.2012.6407894).
- [5] M. Makrehchi, S. Shah, W. Liao, Stock prediction using event-based sentiment analysis, in: *Web Intelligence (WI) and Intelligent Agent Tech-*

- nologies (IAT), 2013 IEEE/WIC/ACM International Joint Conferences on, Vol. 1, 2013, pp. 337–342. doi:10.1109/WI-IAT.2013.48.
- [6] D. Ostrowski, Semantic filtering in social media for trend modeling, in: Semantic Computing (ICSC), 2013 IEEE Seventh International Conference on, 2013, pp. 399–404. doi:10.1109/ICSC.2013.78.
- [7] A. Makazhanov, D. Rafiei, Predicting political preference of twitter users, in: Advances in Social Networks Analysis and Mining (ASONAM), 2013 IEEE/ACM International Conference on, 2013, pp. 298–305.
- [8] H. Hirose, L. Wang, Prediction of infectious disease spread using twitter: A case of influenza, in: Parallel Architectures, Algorithms and Programming (PAAP), 2012 Fifth International Symposium on, 2012, pp. 100–105. doi:10.1109/PAAP.2012.23.
- [9] B. Maitra, State of the retail foreign exchange market, Report (January 2014).
- [10] M. Bank for International Settlements, E. Department, Bis quarterly review - international banking and financial market developments, Report (December 2013).
- [11] P. Papaioannou, L. Russo, G. Papaioannou, C. I. Siettos, Can social microblogging be used to forecast intraday exchange rates?, CoRR abs/1310.5306.
URL <http://arxiv.org/abs/1310.5306>
- [12] M. Makrehchi, S. Shah, W. Liao, Stock prediction using event-based sentiment analysis, in: Web Intelligence (WI) and Intelligent Agent Technologies (IAT), 2013 IEEE/WIC/ACM International Joint Conferences on, Vol. 1, 2013, pp. 337–342. doi:10.1109/WI-IAT.2013.48.
- [13] D. Terrana, A. Augello, G. Pilato, Automatic unsupervised polarity detection on a twitter data stream, in: Semantic Computing (ICSC),

- 2014 IEEE International Conference on, 2014, pp. 128–134. doi:10.1109/ICSC.2014.17.
- [14] N. Godbole, M. Srinivasaiah, S. Skiena, Large-scale sentiment analysis for news and blogs, in: Proceedings of the International Conference on Weblogs and Social Media (ICWSM), 2007.
- [15] S. Fong, S. Deb, I.-W. Chan, P. Vijayakumar, An event driven neural network system for evaluating public moods from online users' comments, in: Applications of Digital Information and Web Technologies (ICADIWT), 2014 Fifth International Conference on the, 2014, pp. 239–243. doi:10.1109/ICADIWT.2014.6814688.
- [16] F. E. Committee, Foreign exchange committee fx volume survey results, Report (July 2014).
- [17] D. Aronson, Evidence-Based Technical Analysis: Applying the Scientific Method and Statistical Inference to Trading Signals, Wiley Trading, Wiley, 2011.
- [18] D. Mocanu, A. Baronchelli, N. Perra, B. Gonçalves, R. Vespignani, The twitter of babel: Mapping world languages through microblogging platforms (2012).
- [19] A. Go, R. Bhayani, L. Huang, Twitter sentiment classification using distant supervision, Tech. rep., Stanford University.
URL <https://sites.google.com/site/twittersentimenthelp/home>
- [20] T. Mikolov, K. Chen, G. Corrado, J. Dean, Efficient estimation of word representations in vector space, CoRR abs/1301.3781.
URL <http://arxiv.org/abs/1301.3781>
- [21] T. Mikolov, W. tau Yih, G. Zweig, Linguistic regularities in continuous space word representations, in: Proceedings of the 2013 Conference

- of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT-2013), Association for Computational Linguistics, 2013.
URL <http://research.microsoft.com/apps/pubs/default.aspx?id=189726>
- [22] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, J. Dean, Distributed representations of words and phrases and their compositionality, in: C. Burges, L. Bottou, M. Welling, Z. Ghahramani, K. Weinberger (Eds.), *Advances in Neural Information Processing Systems 26*, Curran Associates, Inc., 2013, pp. 3111–3119.
- [23] J. J. Faraway, *Extending the linear model with R: generalized linear, mixed effects and nonparametric regression models*, Vol. 66, CRC press, 2006.
- [24] M. Y. Park, *Generalized linear models with regularization*, Ph.D. thesis, Stanford University. Department of Statistics (2006).
- [25] M. Pal, *Multiclass approaches for support vector machine based land cover classification*, CoRR abs/0802.2411.
URL <http://arxiv.org/abs/0802.2411>
- [26] L. Breiman, Random forests, *Mach. Learn.* 45 (1) (2001) 5–32. doi:10.1023/A:1010933404324.
URL <http://dx.doi.org/10.1023/A:1010933404324>
- [27] S. Džeroski, B. Ženko, Is combining classifiers with stacking better than selecting the best one?, *Mach. Learn.* 54 (3) (2004) 255–273. doi:10.1023/B:MACH.0000015881.36452.6e.
URL <http://dx.doi.org/10.1023/B:MACH.0000015881.36452.6e>
- [28] T. Wilson, J. Wiebe, P. Hoffmann, Recognizing contextual polarity in phrase-level sentiment analysis, in: *Proceedings of the Conference on*

Human Language Technology and Empirical Methods in Natural Language Processing, HLT '05, Association for Computational Linguistics, Stroudsburg, PA, USA, 2005, pp. 347–354. doi:10.3115/1220575.1220619.

URL <http://dx.doi.org/10.3115/1220575.1220619>

- [29] AMF France, Study of investment performance of individuals trading in cfd's and forex in france (2014).

Dodatek A

Seznam rezultatov simulacij

Seznam uporabljenih metod za predstavitev besedil in parametrov, razvrščen glede na stanje sredstev po zaključku simulacije trgovanja v padajočem vrstnem redu. Simulacije smo izvedli v opisanem simulacijskem okolju, ki posnema resnično trgovanje v obdobju štirih let in osmih mesecev. Seznam je omejen na najboljših 600 eksperimentov.

metoda	uč. primerov	sredstva (USD)	poslov	ACC	razred	uteži	<i>penalty</i>	<i>alpha</i>	
vreča nizov	100	2378.02	468	0.63	atr	w10	12	1	
vreča nizov	120	2375.74	461	0.63	atr	w10	12	1	
vreča nizov	120	2374.19	510	0.63	atr	w9	12	1	
uvrščanje word2vec	60	2365.65	564	0.63	atr	w10	12	1	
povprečje word2vec	140	2359.45	324	0.60	atr	w10	12	1	
vreča nizov	80	2355.92	503	0.65	atr	w9	12	1	
vreča nizov	60	2353.35	499	0.64	atr	w9	12	1	
vreča nizov	80	2349.47	472	0.63	atr	w10	12	1	
vreča nizov	140	2332.85	462	0.63	atr	w10	12	1	
vreča besed	60	2330.67	524	0.64		w4	12	1	
vreča besed	140	2329.69	472	0.64		w10	12	1	
vreča besed	120	2329.54	474	0.64		w10	12	1	
povprečje word2vec	120	2327.27	346	0.59	atr	w9	12	1	
vreča nizov	80	2312.84	478	0.64	atr	w9	*	12	1
vreča besed	80	2307.85	480	0.65		w10	12	1	
vreča nizov	140	2304.92	475	0.64	atr	w9	*	12	1
vreča nizov	120	2303.06	479	0.63	atr	w9	*	12	1
uvrščanje word2vec	100	2291.95	503	0.61	atr	w10	12	1	
povprečje word2vec	80	2282.35	357	0.60	atr	w10	12	1	
vreča nizov	140	2281.47	496	0.64	atr	w9	12	1	
vreča besed	60	2271.10	643	0.64		w4	11	1e-6	
vreča besed	120	2260.63	516	0.65		w9	12	1	
vreča besed	60	2255.21	553	0.65		w3	12	1	

znak * pomeni, da smo uporabili tehniko sprotnega izbora atributov

vreča nizov	100	2248.96	482	0.64	atr	w9	*	12	1
vreča besed	60	2236.42	482	0.64		w10		12	1
povprečje word2vec	60	2235.49	364	0.60	atr	w10		12	1
vreča besed	100	2232.55	480	0.64		w10		12	1
vreča nizov	100	2232.32	506	0.64	atr	w9		12	1
vreča nizov	120	2230.32	452	0.63	atr	w10	*	12	1
vreča besed	140	2210.03	510	0.64		w9		12	1
povprečje word2vec	140	2202.57	332	0.60	atr	w9		12	1
vreča nizov	180	2186.75	531	0.62	atr	w3		12	1
vreča nizov	140	2186.12	450	0.63	atr	w10	*	12	1
vreča nizov	80	2183.37	468	0.63	atr	w10	*	12	1
vreča besed	120	2180.51	550	0.66		w3		12	1
uvrščanje word2vec	80	2178.05	499	0.63	atr	w10		12	1
vreča nizov	60	2159.89	484	0.63	atr	w10		12	1
vreča besed	80	2149.09	563	0.65		w3		12	1
uvrščanje word2vec	140	2148.65	407	0.59	atr	w10		12	1
povprečje word2vec	100	2148.36	354	0.59	atr	w10		12	1
vreča nizov	180	2146.91	469	0.63	atr	w9	*	12	1
vreča besed	80	2142.69	581	0.63		w9		11	1e-6
vreča besed	140	2137.99	547	0.65		w3		12	1
vreča besed	100	2134.33	561	0.65		w3		12	1
vreča nizov	60	2134.29	551	0.63	atr	w3		12	1
povprečje word2vec	120	2134.10	349	0.58	atr	w10		12	1
vreča nizov	180	2130.45	443	0.63	atr	w10	*	12	1
povprečje word2vec	180	2127.56	317	0.60	atr	w10		12	1
vreča besed	80	2119.77	535	0.64		w9		12	1
vreča besed	80	2119.37	625	0.64		w4		11	1e-6
vreča nizov	100	2119.03	469	0.63	atr	w10	*	12	1
vreča nizov	80	2112.05	503	0.61	atr	w3	*	12	1
vreča besed	80	2111.17	634	0.64		w5		11	1e-6
vreča besed	100	2107.82	527	0.64		w9		12	1
vreča nizov	140	2107.29	533	0.63	atr	w3		12	1
vreča besed	120	2099.17	621	0.63		w4		11	1e-6
vreča nizov	60	2054.05	496	0.62	atr	w3	*	12	1
vreča nizov	120	2051.73	136	0.56	atr	w8	*	12	1
vreča nizov	100	2051.10	543	0.63	atr	w3		12	1
vreča besed	80	2050.87	532	0.65		w4		12	1
vreča nizov	60	2046.42	492	0.63	atr	w9	*	12	1
vreča besed	180	2039.34	545	0.65		w3		12	1
vreča besed	180	2034.34	495	0.64		w9		12	1
vreča besed	180	2025.69	588	0.62		w5		11	1e-6
vreča besed	140	2025.02	544	0.64		w9		11	1e-6
vreča nizov	60	2015.62	469	0.64	atr	w10	*	12	1
vreča nizov	140	2014.41	505	0.61	atr	w4		12	1
vreča nizov	100	2005.12	597	0.57	atr	w10		11	1e-6
vreča nizov	60	1999.12	501	0.61	atr	w4		12	1
povprečje word2vec	100	1996.67	354	0.60	atr	w9		12	1
vreča nizov	140	1991.62	107	0.56	atr	w8	*	12	1
vreča nizov	180	1988.23	393	0.57	atr	w5	*	12	1
vreča besed	100	1982.66	536	0.65		w4		12	1
vreča nizov	120	1979.79	501	0.61	atr	w4		12	1
vreča besed	120	1973	531	0.65		w4		12	1
vreča besed	60	1970.13	537	0.64		w9		12	1
vreča besed	120	1957.63	220	0.57		w6		12	1
vreča nizov	180	1955.94	471	0.60	atr	w4	*	12	1
vreča nizov	120	1948.78	500	0.61	atr	w3	*	12	1
uvrščanje word2vec	140	1945.73	410	0.60	atr	w9		12	1
vreča nizov	120	1942.59	536	0.62	atr	w3		12	1
uvrščanje word2vec	100	1941.07	495	0.60	atr	w9		12	1
povprečje word2vec	180	1939.24	204	0.54	atr	w2		12	1
vreča besed	60	1939.08	586	0.64		w9		11	1e-6
povprečje word2vec	180	1938.04	318	0.59	atr	w9		12	1
vreča nizov	80	1933.27	569	0.59	atr	w10		11	1e-6
vreča nizov	80	1930.79	546	0.62	atr	w3		12	1

vreča besed	180	1925.19	540	0.63		w2		l2	1
vreča nizov	140	1922.29	599	0.59	atr	w9		l1	1e-6
povprečje word2vec	80	1921.03	226	0.41	change	w3		l2	1
vreča besed	180	1918.04	538	0.64		w9		l1	1e-6
povprečje word2vec	120	1918.01	74	0.55	atr	w8		l2	1
vreča nizov	140	1917.90	200	0.56	atr	w7		l2	1
vreča nizov	180	1917.65	496	0.62	atr	w3	*	l2	1
vreča besed	120	1913.69	189	0.56		w7		l2	1
vreča besed	140	1911.03	525	0.64		w4		l2	1
vreča nizov	100	1907.84	506	0.61	atr	w4		l2	1
vreča besed	140	1905.28	627	0.65		w4		l1	1e-6
vreča nizov	120	1902.94	442	0.59	atr	w5		l2	1
povprečje word2vec	60	1897.69	355	0.59	atr	w9		l2	1
vreča besed	140	1894.23	226	0.56		w6		l2	1
vreča nizov	100	1889.82	514	0.61	atr	w3	*	l2	1
vreča nizov	80	1887.41	459	0.60	atr	w4	*	l2	1
vreča besed	120	1884.71	439	0.61		w5		l2	1
vreča besed	120	1880.52	157	0.57		w8		l2	1
vreča nizov	60	1874.12	602	0.60	atr	w9		l1	1e-6
vreča besed	100	1871.54	459	0.61		w5		l2	1
uvrščanje word2vec	80	1863.37	526	0.60	atr	w9		l2	1
vreča nizov	80	1862.67	573	0.61	atr	w9		l1	1e-6
vreča nizov	140	1861.18	144	0.56	atr	w7	*	l2	1
vreča besed	140	1855.21	622	0.62		w5		l1	1e-6
vreča nizov	80	1849.82	651	0.60	atr	w5		l1	1e-6
vreča besed	140	1849.45	562	0.63		w10		l1	1e-6
uvrščanje word2vec	60	1841.06	578	0.63	atr	w9		l2	1
vreča nizov	140	1833.20	506	0.60	atr	w3	*	l2	1
vreča nizov	120	1828.82	438	0.61	atr	w4	*	l2	1
vreča nizov	80	1827.86	506	0.61	atr	w4		l2	1
vreča nizov	120	1825.05	503	0.60	atr	w2		l2	1
povprečje word2vec	80	1816.35	203	0.42	change	w2		l2	1
vreča nizov	120	1816.23	192	0.56	atr	w7		l2	1
povprečje word2vec	80	1815.17	348	0.59	atr	w9		l2	1
povprečje word2vec	120	1813.09	82	0.55	atr	w7		l2	1
vreča besed	140	1811.65	213	0.56		w7		l2	1
vreča nizov	120	1811.56	178	0.56	atr	w8		l2	1
uvrščanje word2vec	120	1803.76	466	0.57	atr	w9		l2	1
vreča nizov	180	1803.24	152	0.56	atr	w7	*	l2	1
povprečje word2vec	140	1802.82	214	0.54	atr	w2		l2	1
povprečje word2vec	120	1800.08	188	0.57	atr	w4		l2	1
uvrščanje word2vec	180	1799.32	392	0.60	atr	w9		l2	1
vreča besed	80	1799.21	585	0.61		w10		l1	1e-6
povprečje word2vec	80	1797.14	243	0.55	atr	w3		l2	1
vreča nizov	120	1795.11	164	0.55	atr	w7	*	l2	1
vreča nizov	180	1790.83	137	0.56	atr	w8	*	l2	1
vreča nizov	100	1788.82	615	0.59	atr	w9		l1	1e-6
povprečje word2vec	100	1785.81	221	0.55	atr	w4		l2	1
povprečje word2vec	140	1782.18	197	0.54	atr	w4		l2	1
vreča besed	140	1774.65	442	0.61		w5		l2	1
vreča nizov	140	1767.87	658	0.59	atr	w5		l1	1e-6
vreča nizov	80	1766.92	370	0.57	atr	w5	*	l2	1
vreča besed	80	1766.76	521	0.63		w2		l2	1
povprečje word2vec	120	1760.28	219	0.55	atr	w2		l2	1
vreča nizov	140	1755.66	184	0.55	atr	w8		l2	1
vreča nizov	120	1752.98	267	0.55	atr	w6		l2	1
vreča besed	120	1752.49	559	0.65		w9		l1	1e-6
uvrščanje word2vec	120	1749.77	507	0.58	atr	w10		l2	1
vreča nizov	120	1747.60	608	0.60	atr	w9		l1	1e-6
vreča nizov	80	1744.95	655	0.62	atr	w4		l1	1e-6
vreča nizov	180	1743.59	658	0.58	atr	w5		l1	1e-6
vreča besed	100	1742.80	578	0.62		w10		l1	1e-6
vreča besed	100	1742.36	535	0.63		w2		l2	1
povprečje word2vec	140	1739.95	133	0.53	atr	w5		l2	1

povprečje word2vec	80	1726.73	176	0.42	change	w4	12	1	
vreča nizov	180	1726.11	498	0.59	atr	w2	12	1	
povprečje word2vec	100	1726	189	0.41	change	w4	12	1	
vreča besed	100	1724.71	637	0.62		w5	11	1e-6	
vreča besed	100	1716.47	608	0.62		w9	11	1e-6	
vreča nizov	100	1713.17	682	0.60	atr	w4	11	1e-6	
povprečje word2vec	140	1712.35	241	0.41	change	w4	12	1	
vreča besed	100	1710.47	631	0.64		w4	11	1e-6	
vreča nizov	180	1708.87	644	0.60	atr	w4	11	1e-6	
povprečje word2vec	80	1706.85	144	0.41	change	w5	12	1	
vreča nizov	80	1696.21	421	0.58	atr	w5	12	1	
povprečje word2vec	120	1695.98	294	0.40	change	w3	12	1	
povprečje word2vec	120	1695.73	130	0.55	atr	w5	12	1	
povprečje word2vec	120	1682.11	91	0.54	atr	w6	12	1	
uvrščanje word2vec	120	1677.09	483	0.57	atr	w7	12	1	
povprečje word2vec	120	1671.73	134	0.41	change	w5	12	1	
povprečje word2vec	120	1671.61	225	0.41	change	w4	12	1	
vreča nizov	120	1668.08	644	0.60	atr	w4	11	1e-6	
vreča nizov	120	1663.92	237	0.55	atr	w6	*	12	1
uvrščanje word2vec	80	1654.89	592	0.54	atr	w2	12	1	
vreča nizov	140	1650.55	446	0.56	atr	w5	12	1	
povprečje word2vec	100	1646.46	127	0.40	change	w5	12	1	
povprečje word2vec	120	1639.48	92	0.41	change	w7	12	1	
vreča nizov	180	1629.92	480	0.59	atr	w2	*	12	1
vreča nizov	180	1626.12	574	0.60	atr	w9	11	1e-6	
vreča nizov	140	1624.21	281	0.54	atr	w6	12	1	
povprečje word2vec	140	1624.16	317	0.39	change	w3	12	1	
vreča besed	180	1623.85	566	0.61		w10	11	1e-6	
vreča besed	60	1618.57	508	0.61		w2	12	1	
vreča besed	120	1615.24	269	0.55		w1	12	1	
povprečje word2vec	120	1611.68	94	0.41	change	w8	12	1	
povprečje word2vec	120	1609.72	87	0.53	atr	w1	12	1	
vreča besed	120	1604.17	533	0.62		w2	12	1	
vreča nizov	180	1604.05	584	0.58	atr	w10	11	1e-6	
vreča nizov	140	1599.32	523	0.60	atr	w2	12	1	
povprečje word2vec	120	1593.23	72	0.41	change	w6	12	1	
povprečje word2vec	80	1587.13	199	0.54	atr	w4	12	1	
vreča nizov	80	1585.73	432	0.58	atr	w2	*	12	1
vreča nizov	80	1583.92	505	0.60	atr	w2	12	1	
vreča besed	180	1580.90	630	0.62		w4	11	1e-6	
vreča besed	120	1578.83	666	0.61		w5	11	1e-6	
vreča besed	60	1576.95	651	0.61		w5	11	1e-6	
tehnični1	60	1576.62	682	0.34	change	w10	12	1	
vreča nizov	60	1569.85	466	0.59	atr	w4	*	12	1
vreča besed	80	1569.69	457	0.60		w5	12	1	
vreča nizov	60	1567.37	664	0.60	atr	w5	11	1e-6	
povprečje word2vec	80	1558.70	127	0.41	change	w8	12	1	
vreča besed	140	1558.18	209	0.55		w8	12	1	
vreča besed	120	1556.15	561	0.62		w10	11	1e-6	
vreča besed	60	1554.26	440	0.59		w5	12	1	
povprečje word2vec	100	1553.65	280	0.55	atr	w3	12	1	
vreča nizov	140	1552.69	643	0.59	atr	w4	11	1e-6	
povprečje word2vec	140	1550.11	71	0.54	atr	w6	12	1	
povprečje word2vec	140	1546.23	71	0.54	atr	w8	12	1	
vreča nizov	60	1543.21	415	0.56	atr	w5	12	1	
povprečje word2vec	100	1542.07	262	0.42	change	w2	12	1	
vreča nizov	180	1540.51	268	0.55	atr	w6	*	12	1
tehnični1	80	1533.65	697	0.38	change	w2	12	1	
vreča nizov	120	1532.60	666	0.59	atr	w5	11	1e-6	
povprečje word2vec	180	1531.43	231	0.40	change	w4	12	1	
povprečje word2vec	80	1531.40	138	0.40	change	w1	12	1	
povprečje word2vec	100	1530.25	210	0.53	atr	w2	12	1	
vreča nizov	60	1530.02	212	0.54	atr	w8	12	1	
povprečje word2vec	80	1524.74	180	0.52	atr	w5	12	1	

povprečje word2vec	140	1522.40	248	0.40	change	w2		12	1
povprečje word2vec	80	1521.53	138	0.40	change	w7		12	1
povprečje word2vec	120	1519.99	210	0.41	change	w2		12	1
vreča besed	180	1519.15	325	0.53		w1		12	1
uvrščanje word2vec	60	1511.51	636	0.55	atr	w2		12	1
vreča besed	60	1509.76	597	0.61		w10		11	1e-6
uvrščanje word2vec	100	1508.61	528	0.56	atr	w4		12	1
vreča nizov	140	1506.38	278	0.54	atr	w6	*	12	1
vreča nizov	140	1503.36	461	0.59	atr	w4	*	12	1
povprečje word2vec	100	1500.56	152	0.54	atr	w5		12	1
uvrščanje word2vec	120	1499.09	529	0.54	atr	w5		12	1
povprečje word2vec	180	1494.50	232	0.54	atr	w4		12	1
vreča nizov	60	1491.80	622	0.59	atr	w10		11	1e-6
uvrščanje word2vec	80	1486.34	539	0.55	atr	w3		12	1
povprečje word2vec	100	1485.98	107	0.53	atr	w8		12	1
vreča nizov	180	1474.86	405	0.52	atr	w1		12	1
povprečje word2vec	120	1467.47	246	0.56	atr	w3		12	1
povprečje word2vec	100	1457.69	107	0.41	change	w7		12	1
vreča nizov	60	1455.86	230	0.54	atr	w7		12	1
povprečje word2vec	100	1452.29	112	0.53	atr	w7		12	1
povprečje word2vec	80	1447.55	134	0.41	change	w6		12	1
vreča besed	140	1444.62	536	0.63		w2		12	1
vreča nizov	80	1444.61	231	0.54	atr	w8		12	1
povprečje word2vec	140	1441.05	77	0.40	change	w6		12	1
povprečje word2vec	140	1439.79	86	0.41	change	w1		12	1
vreča nizov	120	1439.12	473	0.60	atr	w2	*	12	1
povprečje word2vec	140	1434.53	75	0.54	atr	w7		12	1
vreča nizov	60	1432.70	233	0.53	atr	w6	*	12	1
povprečje word2vec	180	1432.25	333	0.39	change	w3		12	1
povprečje word2vec	120	1431.04	91	0.41	change	w1		12	1
vreča nizov	120	1429.24	407	0.56	atr	w5	*	12	1
povprečje word2vec	80	1428.62	240	0.53	atr	w2		12	1
povprečje word2vec	100	1421.62	281	0.41	change	w3		12	1
povprečje word2vec	180	1419.02	127	0.52	atr	w5		12	1
vreča nizov	120	1416.25	596	0.58	atr	w10		11	1e-6
povprečje word2vec	140	1416.17	141	0.41	change	w5		12	1
vreča nizov	60	1412.51	216	0.54	atr	w7	*	12	1
povprečje word2vec	140	1406.70	274	0.54	atr	w3		12	1
vreča nizov	60	1403.53	259	0.53	atr	w6		12	1
vreča nizov	100	1403.22	520	0.60	atr	w2		12	1
vreča nizov	140	1402.80	406	0.56	atr	w5	*	12	1
povprečje word2vec	100	1391.16	112	0.41	change	w6		12	1
povprečje word2vec	100	1387.81	128	0.53	atr	w6		12	1
povprečje word2vec	100	1379.17	124	0.53	atr	w1		12	1
vreča nizov	60	1373.18	680	0.61	atr	w4		11	1e-6
vreča nizov	60	1371.38	477	0.58	atr	w2		12	1
povprečje word2vec	60	1370.75	205	0.53	atr	w5		12	1
vreča nizov	80	1370.26	260	0.55	atr	w7		12	1
vreča besed	100	1361.89	214	0.55		w8		12	1
vreča nizov	100	1355.87	481	0.60	atr	w4	*	12	1
uvrščanje word2vec	100	1355.86	614	0.57	atr	w2		12	1
uvrščanje word2vec	120	1353.64	553	0.53	atr	w2		12	1
uvrščanje word2vec	100	1345.80	590	0.55	atr	w5		12	1
uvrščanje word2vec	80	1345.45	581	0.55	atr	w5		12	1
uvrščanje word2vec	120	1341.29	442	0.55	atr	w8		12	1
vreča besed	100	1340.61	240	0.56		w7		12	1
vreča nizov	100	1334.39	464	0.59	atr	w2	*	12	1
vreča nizov	60	1332.25	210	0.53	atr	w8	*	12	1
povprečje word2vec	100	1330.56	118	0.40	change	w8		12	1
vreča besed	140	1330.39	270	0.54		w1		12	1
povprečje word2vec	60	1330.27	301	0.39	change	w3		12	1
vreča nizov	100	1327.77	214	0.55	atr	w8		12	1
tehnični1	80	1326.71	717	0.38	change	w9		12	1
vreča besed	60	1323.11	255	0.53		w7		12	1

uvrščanje word2vec	80	1317.60	462	0.35	change	w10	12	1	
tehnični1	60	1317.60	731	0.37	change	w1	12	1	
tehnični1	80	1317.05	707	0.39	change	w1	12	1	
vreča nizov	100	1307.14	654	0.59	atr	w5	11	1e-6	
vreča besed	100	1298.23	272	0.55		w6	12	1	
povprečje word2vec	140	1295.71	96	0.40	change	w7	12	1	
vreča nizov	180	1293.80	399	0.52	atr	w1	*	12	1
tehnični2	140	1287.87	383	0.48		w9	12	1	
vreča besed	80	1285.78	244	0.54		w8	12	1	
povprečje word2vec	120	1282.54	400	0.39	change	w9	12	1	
vreča nizov	80	1282.35	324	0.51	atr	w1	*	12	1
povprečje word2vec	180	1280.93	230	0.40	change	w2	12	1	
tehnični1	60	1278.41	712	0.37	change	w3	12	1	
povprečje word2vec	180	1273.18	44	0.53	atr	w8	12	1	
uvrščanje word2vec	120	1271.51	510	0.57	atr	w3	12	1	
vreča nizov	120	1269.75	397	0.53	atr	w1	12	1	
povprečje word2vec	60	1267.07	258	0.54	atr	w3	12	1	
uvrščanje word2vec	80	1266.85	562	0.56	atr	w4	12	1	
povprečje word2vec	180	1265.90	384	0.38	change	w9	12	1	
vreča nizov	80	1265.49	215	0.53	atr	w8	*	12	1
uvrščanje word2vec	120	1258.79	643	0.39	change	w1	12	1	
povprečje word2vec	140	1256.44	60	0.53	atr	w1	12	1	
vreča nizov	80	1255.34	338	0.51	atr	w1	12	1	
uvrščanje word2vec	180	1251.16	462	0.54	atr	w7	12	1	
vreča nizov	140	1250.47	465	0.58	atr	w2	*	12	1
vreča nizov	80	1245.26	221	0.54	atr	w7	*	12	1
uvrščanje word2vec	140	1244.53	487	0.54	atr	w4	12	1	
povprečje word2vec	60	1244.06	271	0.53	atr	w2	12	1	
uvrščanje word2vec	180	1241.29	418	0.37	change	w8	12	1	
uvrščanje word2vec	180	1241.09	411	0.56	atr	w8	12	1	
uvrščanje word2vec	180	1241.05	558	0.53	atr	w6	12	1	
uvrščanje word2vec	60	1239.59	508	0.55	atr	w7	12	1	
povprečje word2vec	60	1239.52	254	0.40	change	w4	12	1	
vreča nizov	80	1238.72	358	0.41	change	w6	*	12	1
povprečje word2vec	180	1238.14	164	0.40	change	w5	12	1	
vreča besed	60	1237.25	233	0.54		w8	12	1	
vreča nizov	120	1236.97	574	0.38	change	w4	*	12	1
tehnični1	140	1234.41	685	0.36	change	w4	12	1	
povprečje word2vec	140	1234.11	378	0.39	change	w9	12	1	
vreča nizov	140	1233.61	379	0.53	atr	w1	12	1	
uvrščanje word2vec	60	1232.26	604	0.54	atr	w5	12	1	
tehnični2	60	1231.66	576	0.55		w3	12	1	
uvrščanje word2vec	180	1229.09	535	0.54	atr	w5	12	1	
povprečje word2vec	180	1226.72	299	0.52	atr	w3	12	1	
uvrščanje word2vec	100	1223.93	530	0.38	change	w2	12	1	
povprečje word2vec	180	1220.43	54	0.53	atr	w7	12	1	
povprečje word2vec	60	1220.42	187	0.39	change	w1	12	1	
povprečje word2vec	180	1215.28	112	0.38	change	w1	12	1	
uvrščanje word2vec	80	1214.92	598	0.54	atr	w6	12	1	
vreča nizov	100	1206.82	469	0.58	atr	w5	12	1	
uvrščanje word2vec	60	1205.16	589	0.38	change	w2	12	1	
uvrščanje word2vec	180	1204.75	377	0.55	atr	w10	12	1	
povprečje word2vec	100	1200.97	133	0.40	change	w1	12	1	
uvrščanje word2vec	140	1200.06	513	0.55	atr	w3	12	1	
vreča nizov	120	1200.03	359	0.52	atr	w1	*	12	1
tehnični1	140	1198.49	692	0.37	change	w9	12	1	
vreča nizov	120	1197.59	574	0.37	change	w3	*	12	1
uvrščanje word2vec	180	1193.70	530	0.55	atr	w2	12	1	
tehnični1	100	1188.09	595	0.60	atr	w3	12	1	
povprečje word2vec	180	1187.34	59	0.53	atr	w6	12	1	
tehnični1	120	1186.76	630	0.53	atr	w7	12	1	
uvrščanje word2vec	80	1184.48	787	0.51	atr	w9	12	1e-3	
povprečje word2vec	180	1184.31	102	0.40	change	w6	12	1	
uvrščanje word2vec	120	1181.17	410	0.34	change	w10	12	1	

uvrščanje word2vec	140	1181.16	425	0.33	change	w10		12	1
tehnični1	60	1180.56	575	0.52	atr	w10		12	1
tehnični2	80	1179.97	583	0.56		w3		12	1
povprečje word2vec	140	1178.79	103	0.39	change	w8		12	1
povprečje word2vec	60	1177.59	269	0.39	change	w2		12	1
uvrščanje word2vec	60	1176.62	657	0.51	atr	w6		12	1
vreča nizov	80	1175.86	262	0.54	atr	w6	*	12	1
vreča besed	60	1175.34	267	0.52		w1		12	1
tehnični1	120	1167.81	553	0.54	atr	w8		12	1
tehnični2	100	1166.16	584	0.56		w3		12	1
vreča nizov	100	1158.93	204	0.54	atr	w8	*	12	1
povprečje word2vec	60	1158.70	232	0.53	atr	w4		12	1
vreča nizov	80	1154.82	304	0.54	atr	w6		12	1
vreča nizov	80	1153.94	247	0.40	change	w8	*	12	1
tehnični1	120	1149.82	597	0.52	atr	w10		12	1
vreča nizov	60	1149.49	384	0.54	atr	w5	*	12	1
vreča nizov	120	1145.02	404	0.40	change	w1	*	12	1
vreča nizov	100	1143.30	244	0.55	atr	w7		12	1
vreča nizov	80	1138.65	506	0.39	change	w5	*	12	1
uvrščanje word2vec	180	1137.13	398	0.34	change	w10		12	1
tehnični1	120	1136.02	588	0.60	atr	w3		12	1
tehnični1	120	1133.73	686	0.37	change	w8		12	1
uvrščanje word2vec	100	1133	439	0.56	atr	w8		12	1
povprečje word2vec	180	1131.61	92	0.39	change	w7		12	1
tehnični1	60	1130.74	740	0.36	change	w2		12	1
uvrščanje word2vec	140	1130.46	696	0.49	atr	w1		12	1
povprečje word2vec	180	1128.23	85	0.40	change	w8		12	1
vreča nizov	120	1127.40	595	0.40	change	w2	*	12	1
uvrščanje word2vec	140	1127.26	437	0.55	atr	w8		12	1
uvrščanje word2vec	140	1126.68	659	0.39	change	w1		12	1
uvrščanje word2vec	60	1126.33	506	0.56	atr	w8		12	1
vreča nizov	180	1125.97	505	0.39	change	w5	*	12	1
tehnični1	120	1122.89	647	0.53	atr	w3		12	1
uvrščanje word2vec	180	1121.69	470	0.38	change	w3		12	1
vreča besed	60	1121.51	270	0.54		w6		12	1
uvrščanje word2vec	140	1114.37	589	0.52	atr	w2		12	1
tehnični2	120	1111.70	377	0.47		w9		12	1
vreča nizov	180	1109.34	249	0.39	change	w8	*	12	1
tehnični1	80	1108.97	611	0.51	atr	w10		12	1
povprečje word2vec	80	1108.18	164	0.52	atr	w7		12	1
tehnični1	180	1107.24	674	0.35	change	w3		12	1
vreča besed	80	1106.13	275	0.55		w6		12	1
povprečje word2vec	60	1102.28	181	0.52	atr	w7		12	1
uvrščanje word2vec	180	1099.55	631	0.38	change	w1		12	1
povprečje word2vec	80	1096.36	168	0.51	atr	w6		12	1
vreča nizov	100	1094.99	217	0.55	atr	w7	*	12	1
vreča nizov	180	1088.91	393	0.40	change	w6	*	12	1
uvrščanje word2vec	180	1088.45	521	0.54	atr	w3		12	1
uvrščanje word2vec	140	1080.37	566	0.53	atr	w5		12	1
tehnični1	140	1079.75	682	0.39	change	w8		12	1
tehnični2	120	1079.60	585	0.56		w3		12	1
vreča besed	80	1071.18	274	0.55		w7		12	1
tehnični2	80	1070.04	612	0.60		w4		12	1
vreča nizov	60	1069.46	297	0.52	atr	w1		12	1
tehnični1	80	1066.65	600	0.60	atr	w3		12	1
uvrščanje word2vec	140	1066.40	587	0.52	atr	w6		12	1
tehnični1	140	1061.79	585	0.59	atr	w3		12	1
tehnični1	60	1061.42	750	0.36	change	w9		12	1
povprečje word2vec	180	1059.44	73	0.53	atr	w1		12	1
tehnični1	100	1058.83	694	0.37	change	w3		12	1
povprečje word2vec	100	1055.07	425	0.38	change	w9		12	1
povprečje word2vec	80	1053.27	168	0.51	atr	w8		12	1
tehnični1	140	1051.85	662	0.37	change	w1		12	1
povprečje word2vec	60	1050.05	181	0.52	atr	w8		12	1

tehnični2	60	1049.58	662	0.66		w5		12	1
tehnični1	100	1045.62	684	0.34	change	w2		12	1
vreča nizov	140	1042.28	590	0.41	change	w2	*	12	1
tehnični2	80	1042.27	659	0.66		w5		12	1
povprečje word2vec	60	1040.45	216	0.38	change	w5		12	1
vreča nizov	120	1039.72	537	0.40	change	w5	*	12	1
povprečje word2vec	80	1038.41	419	0.38	change	w9		12	1
uvrščanje word2vec	100	1038.23	518	0.37	change	w4		12	1
tehnični2	80	1037.17	594	0.64		w7		12	1
uvrščanje word2vec	100	1036.42	517	0.56	atr	w3		12	1
vreča nizov	80	1031.22	280	0.40	change	w7	*	12	1
tehnični1	80	1030.97	416	0.56	atr	w9		12	1
uvrščanje word2vec	80	1030.16	469	0.33	change	w9		12	1
tehnični1	80	1027.84	622	0.63	atr	w4		12	1
povprečje word2vec	60	1027.41	776	0.52	atr	w10		11	1e-5
tehnični1	60	1027.04	591	0.51	atr	w9		12	1
tehnični1	180	1026.54	574	0.60	atr	w3		12	1
tehnični1	140	1026.44	375	0.55	atr	w10		12	1
vreča nizov	80	1025.48	586	0.40	change	w2	*	12	1
tehnični2	120	1023	341	0.43		w10		12	1
tehnični1	80	1018.85	660	0.66	atr	w5		12	1
tehnični1	120	1018.17	422	0.56	atr	w9		12	1
tehnični1	100	1015.25	697	0.34	change	w6		12	1
tehnični1	120	1014.33	687	0.52	atr	w6		12	1
uvrščanje word2vec	80	1009.01	527	0.36	change	w5		12	1
vreča nizov	140	1008.35	585	0.59	atr	w10		11	1e-6
tehnični2	60	1007.71	611	0.60		w4		12	1
uvrščanje word2vec	100	1006.79	453	0.34	change	w9		12	1
tehnični1	60	1005.26	590	0.59	atr	w3		12	1
uvrščanje word2vec	180	1004.32	775	0.50	atr	w9		12	1e-3
uvrščanje word2vec	140	1004.13	404	0.36	change	w9		12	1
uvrščanje word2vec	120	1003.98	560	0.55	atr	w4		12	1
tehnični2	60	1003.06	598	0.64		w7		12	1
povprečje word2vec	80	1002.80	167	0.51	atr	w1		12	1
vreča nizov	100	1001.47	605	0.36	change	w3	*	12	1
vreča nizov	100	995.92	611	0.39	change	w2	*	12	1
vreča nizov	140	994.99	395	0.53	atr	w1	*	12	1
vreča nizov	100	993.02	286	0.54	atr	w6		12	1
vreča nizov	100	992.14	420	0.57	atr	w5	*	12	1
tehnični2	100	991.16	601	0.64		w7		12	1
uvrščanje word2vec	80	990.52	564	0.38	change	w2		12	1
vreča nizov	180	990.10	285	0.39	change	w7	*	12	1
tehnični2	80	988.55	541	0.61		w8		12	1
povprečje word2vec	60	987.88	192	0.39	change	w6		12	1
vreča nizov	120	985.26	240	0.41	change	w8	*	12	1
vreča besed	80	984.65	306	0.51		w1		12	1
vreča nizov	60	981.96	297	0.38	change	w7	*	12	1
povprečje word2vec	60	981.50	177	0.39	change	w7		12	1
vreča nizov	100	981.33	576	0.37	change	w4	*	12	1
uvrščanje word2vec	80	979.09	480	0.57	atr	w8		12	1
vreča nizov	180	977.37	473	0.39	change	w1	*	12	1
povprečje word2vec	60	976.12	189	0.51	atr	w6		12	1
tehnični2	180	975.23	573	0.55		w3		12	1
tehnični1	60	975.08	612	0.66	atr	w7		12	1
tehnični1	140	974.96	696	0.35	change	w5		12	1
tehnični1	140	973.09	411	0.57	atr	w9		12	1
povprečje word2vec	180	972.91	409	0.35	change	w10		12	1
tehnični1	100	971.25	605	0.65	atr	w7		12	1
tehnični1	120	970.86	725	0.52	atr	w1		12	1
vreča nizov	60	969.31	270	0.38	change	w8	*	12	1
uvrščanje word2vec	140	967.66	530	0.39	change	w2		12	1
tehnični2	140	967.52	581	0.56		w3		12	1
vreča nizov	140	965.24	498	0.40	change	w5	*	12	1
tehnični2	100	964.50	653	0.66		w5		12	1

tehnični1	120	964.31	711	0.36	change	w2		12	1
tehnični1	140	964.27	625	0.52	atr	w7		12	1
vreča nizov	100	963.57	408	0.52	atr	w1		12	1
tehnični2	60	963.36	418	0.46		w9		12	1
tehnični1	180	962.68	510	0.52	atr	w9		12	1
povprečje word2vec	180	962.53	754	0.54	atr	w9		11	1e-5
uvrščanje word2vec	100	958.25	798	0.51	atr	w9		12	1e-3
tehnični1	100	957.41	618	0.62	atr	w4		12	1
povprečje word2vec	140	954.28	409	0.35	change	w10		12	1
tehnični1	140	952.66	727	0.36	change	w6		12	1
vreča nizov	80	950.82	425	0.40	change	w1	*	12	1
uvrščanje word2vec	100	948.94	499	0.56	atr	w7		12	1
tehnični1	80	948.13	726	0.41	change	w8		12	1
povprečje word2vec	60	946.94	183	0.51	atr	w1		12	1
uvrščanje word2vec	180	945.44	434	0.36	change	w9		12	1
vreča nizov	60	944.69	591	0.35	change	w3	*	12	1
vreča nizov	60	941.11	599	0.38	change	w2	*	12	1
tehnični1	100	940.96	680	0.51	atr	w3		12	1
vreča nizov	140	938.89	469	0.40	change	w1	*	12	1
vreča nizov	60	938.25	580	0.36	change	w4	*	12	1
tehnični2	140	937.16	349	0.42		w10		12	1
uvrščanje word2vec	100	932.64	815	0.48	atr	w10		12	1e-3
vreča nizov	100	932.55	262	0.39	change	w8	*	12	1
tehnični1	60	931.03	754	0.36	change	w6		12	1
tehnični2	100	930.79	613	0.60		w4		12	1
tehnični1	100	929.88	657	0.66	atr	w5		12	1
tehnični2	120	929.59	598	0.64		w7		12	1
vreča nizov	100	928.07	290	0.54	atr	w6	*	12	1
tehnični2	180	927.38	381	0.47		w9		12	1
tehnični2	60	926.95	545	0.61		w8		12	1
tehnični2	120	926.39	607	0.61		w4		12	1
tehnični1	60	925.87	622	0.62	atr	w4		12	1
vreča nizov	60	925.57	578	0.34	change	w9	*	12	1
povprečje word2vec	60	924.05	189	0.39	change	w8		12	1
tehnični1	60	923.49	705	0.40	change	w8		12	1
vreča nizov	140	923.37	563	0.35	change	w3	*	12	1
tehnični1	100	922.27	779	0.70	atr	w1		12	1
vreča nizov	140	919.29	248	0.41	change	w8	*	12	1
tehnični1	120	918.26	709	0.35	change	w3		12	1
tehnični2	100	917.10	781	0.70		w1		12	1
vreča nizov	60	914.28	460	0.56	atr	w2	*	12	1
tehnični1	120	913.59	599	0.65	atr	w7		12	1
tehnični1	60	913.40	669	0.66	atr	w5		12	1
tehnični1	120	912.44	701	0.38	change	w1		12	1
vreča nizov	100	912.16	545	0.38	change	w5	*	12	1
vreča nizov	60	907.05	309	0.51	atr	w1	*	12	1
vreča nizov	140	905.85	303	0.41	change	w7	*	12	1
tehnični1	120	904.35	612	0.62	atr	w4		12	1
uvrščanje word2vec	180	900.75	665	0.52	atr	w1		12	1
vreča nizov	100	898.61	439	0.40	change	w1	*	12	1
tehnični2	120	898.29	645	0.66		w5		12	1
tehnični1	180	895.44	683	0.38	change	w1		12	1
vreča nizov	80	892.76	574	0.37	change	w4	*	12	1
tehnični1	80	890.54	708	0.68	atr	w6		12	1
vreča nizov	140	889.40	413	0.39	change	w6	*	12	1
uvrščanje word2vec	140	888.15	490	0.36	change	w4		12	1
tehnični1	80	887.75	611	0.66	atr	w7		12	1
tehnični1	100	886.58	436	0.54	atr	w10		12	1
tehnični2	80	886.46	790	0.71		w1		12	1
tehnični1	80	886.46	790	0.70	atr	w1		12	1
tehnični2	60	884.55	712	0.68		w6		12	1
tehnični2	120	883.26	774	0.70		w1		12	1
tehnični1	120	882.51	770	0.70	atr	w1		12	1
vreča nizov	120	881.04	562	0.34	change	w9	*	12	1

vreča besed	100	880.78	318	0.52		w1	12	1	
tehnični1	120	880.01	652	0.66	atr	w5	12	1	
uvrščanje word2vec	140	877.98	448	0.35	change	w8	12	1	
uvrščanje word2vec	60	877.76	615	0.55	atr	w3	12	1	
tehnični1	140	876.42	759	0.70	atr	w1	12	1	
vreča nizov	120	872.87	268	0.41	change	w7	*	12	1
tehnični1	100	872.61	596	0.50	atr	w9	12	1	
uvrščanje word2vec	120	872.61	456	0.37	change	w8	12	1	
tehnični2	140	871.25	761	0.70		w1	12	1	
tehnični1	140	870.13	595	0.65	atr	w7	12	1	
uvrščanje word2vec	100	869.27	626	0.53	atr	w6	12	1	
tehnični2	80	868.74	709	0.68		w6	12	1	
tehnični1	60	867.96	713	0.34	change	w5		12	1
vreča nizov	100	867.31	376	0.50	atr	w1	*	12	1
uvrščanje word2vec	60	867.13	495	0.34	change	w10	12	1	
tehnični1	120	865.27	663	0.51	atr	w4	12	1	
tehnični1	100	864.73	599	0.52	atr	w8	12	1	
uvrščanje word2vec	80	863.70	628	0.38	change	w1	12	1	
vreča nizov	140	861.79	558	0.37	change	w4	*	12	1
uvrščanje word2vec	100	860.66	486	0.32	change	w10	12	1	
tehnični2	140	858.50	594	0.63		w7	12	1	
vreča nizov	100	857.45	281	0.39	change	w7	*	12	1
tehnični1	100	856.38	617	0.51	atr	w10	12	1	
tehnični1	180	856.22	413	0.58	atr	w9	12	1	
tehnični1	60	856.12	715	0.68	atr	w6	12	1	
uvrščanje word2vec	180	855.32	500	0.55	atr	w4	12	1	
tehnični1	100	855.10	696	0.51	atr	w4	12	1	
tehnični1	140	853.18	607	0.62	atr	w4	12	1	
uvrščanje word2vec	180	851.47	519	0.36	change	w2	12	1	
tehnični1	60	851.30	799	0.70	atr	w1	12	1	
tehnični1	100	850.82	705	0.36	change	w10	12	1	
tehnični1	120	850.45	392	0.55	atr	w10	12	1	
tehnični2	100	847.23	377	0.43		w10	12	1	
povprečje word2vec	140	845.65	766	0.50	atr	w9	11	1e-5	
tehnični1	60	843.55	755	0.37	change	w7	12	1	
vreča nizov	120	839.35	522	0.32	change	w10	*	12	1
tehnični1	180	839.05	586	0.51	atr	w3	12	1	
vreča nizov	80	835.66	582	0.36	change	w3	*	12	1
povprečje word2vec	100	835.35	440	0.35	change	w10	12	1	
tehnični1	140	828.92	647	0.66	atr	w5	12	1	
tehnični2	140	827.84	647	0.65		w5	12	1	
tehnični1	80	826.42	622	0.51	atr	w9	12	1	
tehnični2	100	822.24	554	0.61		w8	12	1	
vreča nizov	100	820.26	367	0.40	change	w6	*	12	1
tehnični2	140	819.93	605	0.60		w4	12	1	
tehnični1	100	816.70	703	0.68	atr	w6	12	1	
vreča nizov	80	813.10	522	0.33	change	w10	*	12	1
uvrščanje word2vec	180	812.61	486	0.37	change	w5	12	1	
tehnični1	140	812.03	563	0.52	atr	w10	12	1	
uvrščanje word2vec	100	807.77	679	0.38	change	w1	12	1	
vreča nizov	60	805.32	490	0.38	change	w5	*	12	1
vreča nizov	140	804.78	549	0.35	change	w9	*	12	1
vreča nizov	120	804.61	355	0.4	change	w6	*	12	1
tehnični1	120	802.32	707	0.37	change	w9	12	1	
vreča nizov	60	800.95	531	0.33	change	w10	*	12	1
tehnični2	100	800.55	706	0.68		w6	12	1	
tehnični2	100	799.55	415	0.46		w9	12	1	
uvrščanje word2vec	180	798.21	463	0.38	change	w7	12	1	
povprečje word2vec	60	797.7	445	0.37	change	w9	12	1	
tehnični1	120	790.09	679	0.36	change	w6	12	1	
uvrščanje word2vec	80	787.74	531	0.54	atr	w7	12	1	
tehnični1	100	787.71	422	0.55	atr	w9	12	1	
uvrščanje word2vec	60	784.84	525	0.35	change	w9	12	1	